

DATA ALLOCATION
IN A
DISTRIBUTED DATABASE ENVIRONMENT

by

Kimberley Ann Johnson

B. S., Western Illinois University, 1979

A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Computer Science

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1988

Approved by:


Major Professor

1.4 Report Organization

This report is organized into 4 chapters. Chapter 2 provides an overview of the problems encountered in data distribution and reviews the relevant literature in this area. Chapter 3 provides an overview of the project's design. Chapter 4 contains an extensive example to demonstrate the functionality and use of the tool. Chapter 5 gives concluding remarks and suggests possible extensions.

CHAPTER 2

THE DATA DISTRIBUTION PROBLEM

2.1 Objective

The objective of this chapter is to summarize the many problems addressed in the literature concerning data distribution. A high level understanding of these problems and their interdependencies will familiarize the reader with the complexity these problems introduce into the distributed database design, even when addressed individually. A great deal of research has been done to address different parts of the data distribution problem. Some of the significant work in these areas will be reviewed. To date however, there has been no overall solution that successfully addresses all of the problems under one combined methodology. Further research continues to look for optimal solutions to these issues.

2.2 Why Distribute ?

In a distributed environment there are clear benefits to be derived in distributing data. One obvious advantage is the ability to store the data at the location(s) where it is most frequently used. This will achieve a faster response time and reduce communication costs in a query intensive application. A second advantage in distributing data is the potential to store what might be a very large database on smaller machines through partitioning or splitting the data (partitioning will be reviewed in detail in a later section). As a final example, the distribution of data increases the overall reliability of the system by the simple fact that all data is not stored on one machine (i.e., subject to single site failure) [ROTH81]. Each of these points demonstrate clear advantages, but maximizing these benefits requires important tradeoffs.

2.3 Distribution Schemes and Associated Problems

Consider the optimal allocation of a file in a distributed environment. Intuition would tell you to

store the data where it is most frequently used. If each site owned and used the data exclusively (i.e., no other site needed any of this data) there would be no design problem. All data updates and queries at this node could be handled locally. In reality though, this is not the case as applications have data that is shared among several users. These applications may also have strict requirements as to response time, reliability and consistency of the data they require. For this reason, careful analysis of the distribution scheme to be used is required.

Distribution schemes generally fall into two main categories: partitioned systems and replicated systems [DRAF80].

2.3.1 Partitioned Systems

Partitioning is the process of "assigning a local object (relation) from the logical schema of the database to several physical objects (files) stored in the database" [BRAY81]. A pure partitioned system has no duplication or replication of data items. Vertical partitioning (or partitioning by structure) divides the data by columns or attributes. Application of vertical partitioning would be desirable in cases where only certain attributes of the logical record are needed at locations. An example taken from [DRAF80] will help illustrate this point. This example considers a relation for orders of parts as follows:

ORDERS (CUST #,CNAME,PART NO,PART DESC,QUANTITY ORDERED)

One site may only be concerned with the PART_NO and PART_DESC, while another site may maintain the CUST_NO and CUST_NAME. In this instance there are benefits derived from partitioning the data vertically, as each user can locally control the information with which they are directly concerned.

Horizontal partitioning divides a relation by occurrence or tuple. This type of partitioning is valuable in cases where files can be distributed based upon given data values. Again, using the

ORDERS relation, each site may need all information in the relation but may only deal with one type of part. In this case all tuples associated with a given part are assigned to that particular node.

2.3.2 Replicated Schemes

Replication is the allocation of a single file of the database to multiple sites. If there is no replication the distributed database problem is significantly reduced in the area of concurrency control and synchronization, but the cost of doing a transaction may increase significantly when the transaction must access data from a number of sites. The opposite extreme is full redundancy where each file is present at every node. This method will optimize response time in the case of queries, but will severely impact update performance and cost as all updates must be propagated to every site. Partial redundancy is the balance between these two extremes. Redundancy is important in achieving much of the promise of distributed systems. Without redundancy reliability goals can only partially be met as the unavailability of even a single file may be seen as "total failure" to some applications. Without redundancy the choice of where to store data is an all or nothing situation. Finally, in reality updates tend to be small and simple where queries tend to be quite complex and involve large amounts of data [ROTH81]. In this type of situation redundancy would have clear benefit but again would increase the complexity of the system from a concurrency and synchronization point of view.

2.3.3 Summary

There are several generalizations that have been made concerning the file allocation and replication issues.

Champine reviews the problem in regard to the size of the file and the percentage of exception rate (or remote request). Figure 2.1 summarizes his view [CHAM81]:

Exception Rate	File Size	Distribution Scheme
—	small	replicate
small	large	partition
large	large	centralize

Distribution Scheme
Figure 2.1

His view simply states that if the file size is small, replicate the file and propagate updates. If the exception rate is small but the file size is large, the best solution is to partition the data in order to place the fragments where they are most frequently used. Finally, if the file size is large and exhibits a high exception rate it may be necessary to centralize the file. In other cases where the size of the file has not been considered, studies have indicated that the maximum number of copies of a file should be one unless the ratio of queries to updates is greater than or equal to 50% (ratio = .5) [MURO85], [CASE72]. Although these views are gross simplifications of the data distribution problem, they appeal to intuition when considering the trade-offs between update costs, communication costs and response time. In addition, they serve to highlight the importance of understanding an application's use of data. It should be clear from the discussion thus far that one cannot reasonably decide where to put data unless they know where it will be accessed from and how often.

In summary, the data distribution problem involves a clear understanding of an application and its use of data. Data allocation however is only one aspect in the design of a distributed processing system. The allocation of hardware and software must be considered in the overall design, along with the allocation of data. The number of design factors increases substantially as decisions must

include such things as network topology, channel bandwidths, number of processors, storage capacities, program locations as well as data locations [HEVN84]. All of these decisions are highly interdependent and involve trade-offs in the areas of reliability, performance, development complexity (synchronization and recovery), growth and overall system cost [MAR84]. Due to the complexity of modeling all of these factors concurrently, most research has considered only individual design problems or the combination of only a few [HEVN84].

The remainder of this chapter will first review the research done in the area of "pure" file allocation and then survey extensions to these earlier models, which take into account hardware considerations and file interdependencies.

2.4 File Allocation Solutions

The first attempts to deal with the problem of optimal file placement used mathematical programming techniques. Nearly all of these models were linear integer programming problems to find optimal solutions. Most of these models were driven by system requirements as to performance, cost, minimum access delays, data management overhead and storage. Specific data models were assumed, with variants applied to different system resources. Generally however, these models worked under the following requirements [CHAM81]:

Given: A description of user demand for service stated as volume of requests from each node of a network to each file

Given: A description of resources available, such as network topology, link capacity, cost of storage, communications cost, etc.

Determine: An assignment of files to nodes which minimizes total costs.

2.4.1 'Pure' File Allocation Models

The simplest file allocation model was introduced by Casey [CASE72]. His model looks for the optimal node assignment of a single file which minimizes the total communications cost, under the assumption of a fully connected network with no response time or memory restrictions imposed.

One important contribution of his model is the distinction made between update and query requests. In the case of queries it is assumed that the copy of the file which minimizes total communication costs is chosen, whereas updates are propagated to all file copies. The model represents the total communication cost as a sum of the cost over individual nodes that result from a given file allocation. Casey's cost equation is paraphrased below [CASE72]:

$$C(I) = \sum_{i \in I} \left[\sum_{j \neq i} (UT)_j (UC)_{jk} + QT_i \min_{j \neq i} QC_{jk} \right] + \sum_{k \in I} STR_k$$

where:

- I = index set of system nodes
- j = index for nodes
- k = index for file locations
- STR = fixed cost of storage for locating files at kth node
- QT = query traffic (emanating from node j)
- UT = update traffic (emanating from node j)
- QC = cost of unit of communication from node j to k for query
- UC = cost of unit of communication from node j

His model demonstrated that query costs decreased as the number of copies of the file were increased; however, a penalty is paid for storage and update costs. Rephrased, if storage costs were low and there were no updates, complete duplication would be cost effective. If storage costs were high and update activity was high, then one copy of the file would be optimal. Casey analyzes this trade-off by examining the cost function as the number of file copies is increased. This is done using a directed graph (referenced as "cost graph") where each vertex is a file assignment and has an associated value from the cost function. The edges of the graph are paths corresponding to the addition of a single file. Casey demonstrates the monotonicity of this graph, implying that it is "sufficient to follow every path of the cost graph until the cost increases, and no more" [CASE72]. This bounds the number of computations in both breadth and depth (commonly referred to as a Branch and Bound Search). Therefore, only a subset of the original tree (of 2 to the power of "n" nodes) needs to be tested. Even with this property, Casey's model

was proven by Eswaran to be NP-complete suggesting that heuristics were needed to efficiently deal with the problem [HEVN84].

Another approach to the problem was introduced by Chu. His model sought to minimize overall operating costs by determining the optimal placement of files under the constraints of response time and storage capacity. No distinction was made between query and update requests [CHU69]. This model again introduced a very large number of variables for even small problems, making it very costly, and in large problems computationally infeasible [CER183].

Several comprehensive reviews of other models which address the pure file allocation problem can be found in [CER183a], [HEVN84], [LEVI79]. Most of the proposed solutions apply different heuristic measures to reduce the computational complexity of the models. In general all of the models in this category assume a completely connected network topology and a completely defined distribution scheme in terms of storage capacity, storage costs, communication costs, frequencies, and user requests [HEVN84]. With all of these parameters well defined they attempt to obtain an optimal design for a very specific problem, namely the allocation of a single file. One model which was introduced by Morgan and Levin deserves special note, as it points out another weakness in the models discussed so far.

Morgan and Levin's model distinguishes itself from others by considering both program and data allocation. Their work points out the importance of considering the dependencies between programs and data in a heterogeneous environment. The dependency points out that while data can move easily from node to node, programs cannot, as in a heterogeneous environment different hardware and system software will exist. In this environment, program execution is limited to certain nodes. This is important as a transaction at one node may invoke a program at another node, which in turn needs data from a third node. Their model considers this restriction, by analyzing the optimal placement of both data and programs. In addition their model also

considers dynamic behavior and uncertain demand. Most of the previous models assumed "static" behavior, meaning that once the frequency of requests was determined they would not change. Also assumed was the complete availability and accuracy of access patterns. Morgan and Levin's model deals with dynamic behavior and incomplete information by minimizing costs over a number of different time periods during which allocation may change [LEVI79]. The practicality of changing the distribution over time has been debated however. As pointed out by [GROSS80], once a distribution scheme is set up it remains fairly static in practice as redistribution could involve a tremendous amount of effort.

Ceri, et. al. and Hevner review similar work done by Fisher and Hochbaum [CERI83a],[HEVN84]. This model is an optimization algorithm for placing multiple copies of programs and databases over a network. They improve the work done by Morgan and Levin by developing several heuristics to generate feasible solutions to the problem and report practical experience.

2.4.2 File and Hardware Allocation Models

Various extensions have been made to the file allocation models discussed above, which explore different sets of assumptions and problems. These approaches relaxed some of the restrictions of other models by addressing such things as channel capacities and network topologies.

Mahmoud and Riordan consider the combined problem of optimal file allocation and channel capacity determination given a fixed network topology. The objective of the model is to minimize communication cost and storage cost subject to network delay and average file availability. The model was a nonlinear integer programming problem making it quite expensive; thus, an efficient heuristic was developed [MAHM76]. Their results have produced reasonably good allocation solutions [HEVN84].

Irani combines the file allocation, network topology and channel capacity allocation into a single problem. The model minimizes the total cost of file storage and communication capacities over different channels. The constraints imposed on the model include a maximum communication delay, minimum availability of single files and a minimum level of network reliability [IRAN79].

As a final example, Hevner reviews the work done by Casey [HEVN84]. Casey extended his original model by including optimal selections of network topology and channel capacities. Due to the size of the problem the topology considered is restricted to tree networks. Again, Casey developed heuristic techniques to solve the problem as the original model was nonlinear and contained integer and continuous variables.

2.4.3 Summary

In summary, earlier models developed in this category focused solely on the file allocation problem, leaving all other factors invariant. Extensions were made to these models to deal with program allocation and hardware allocation design choices. Several generalizations can be applied to these models. Most of these solutions were integer linear programming problems (i.e., each design parameter had a 0 or 1 allocation). These were far too costly, and hence infeasible to run for systems of any real size. Heuristic techniques were therefore developed to deal with the problems. These heuristics sacrifice optimality for practicality of use. Even with the use of these heuristics, the complexity of the problem has limited most research to dealing with the optimization of at most two resources at once [HEVN84].

Other major drawbacks of these models include their assumptions in regard to file usage and partitioning. These models do not address interdependencies between files which appear in realistic databases. They assume access of a given file from a given node, and do not reflect the demand for data access involving more than one file. As pointed out by [ROTH81], consider a join where each of two nodes has a file used in the new relation. A join request involving these

two relations could involve substantial communication costs. A file allocation scheme that places both relations involved in this query at a single site may be far more advantageous than if distributed (Although true in principal, this view may be slightly exaggerated as the amount of data involved could be reduced by a query optimization technique that utilizes semi-joins [HEVN84]). Finally, these models assume that complete files should be the unit of assignment of data to nodes. There is no consideration given to partitioning the files in order to reduce access and storage costs. While these models are important, they often obtain the optimal design for a very specific problem. The use of these models may be highly advantageous after the partitioning problem has been addressed and the physical distributed system is designed.

2.5 File Dependency Methodologies

As stated above, one of the major problems with the pure file allocation models is the assumption of a single file being the unit of distribution with no consideration given to file dependencies. Research in this area proposes methodologies and solution methods that consider the entire database schema as opposed to individual files. Some of the research presents guidelines and classification schemes to be followed during logical distribution design (non-automated). A notable contribution in this area is the work done by Baker [BAKE]. His methodology points out the major issues that must be addressed during logical design. Other contributions propose theoretical models and heuristic algorithms to determine the logical distribution. Examples from each category are reviewed below.

2.5.1 Baker's Model

Baker proposes a methodology in which logical distribution is defined as a "partition of a collection of related applications and their data into a maximum number of groups that have a specified low level of interdependence". 'Intermodal dependencies' arise when an application transaction (or program) requires data from a remote node. The goal then is to minimize the

intermodal dependencies and move toward 'nodal autonomy' or 'nodal interdependency' [BAKE].

In Baker's model, transaction or intergroup dependencies have an orientation, meaning an ownership of the transaction. Several quantitative measures are outlined that can be used in different combinations to classify intergroup dependencies. These include frequency of use, pattern of usage, required currency, level of consistency and timing required and the degree of remote data needed. In addition other classifications can be used to rate the dependencies, such as read versus update requests (where read would have a weaker dependency) or deferrability of the transaction.

Baker's approach to the design is iterative in nature and follows through six steps. These steps and their associated activities can be summarized as follows:

1. *Data Gathering:* Determine the number and types of databases, database structures, relations among the databases and information regarding the application make-up (i.e., split geographically or by function).
2. *Define Application Groupings:* Define application group structure such as order entry, production planning, etc.
3. *Assign Applications to Groupings:* Here a complete application which consists of a set of application programs is assigned to one group only.
4. *Assign Databases to Application Groupings:* Based upon knowledge of the applications and the data they require, assign databases to each application group (i.e., applications that make the most updates to data).
5. *Assign Transactions to Databases:* Here step 4 is ignored, and each transaction is assigned to the group that contains the data that are most closely related to the program or transaction. This is included to overcome the possibility of being assigned to an application group where the data the program most frequently uses is in another group.
6. *Analyze Dependencies and Evaluate Distribution:* The objectives of this step are to minimize the communication traffic and minimize the amount of data that must be copied between nodes. This step will be further described below.

In analyzing dependencies and evaluating the distribution (step 6), the dependency of each pair of groups is calculated from the set of transaction dependencies that occur between groups. These dependencies are characterized on a single transaction type by a) an orientation which relates the local to remote group; (b) the active component, or number which gives the frequency of use of

the transaction per day and (c) the passive component, or number which represents the number of bytes of remote data accessed by the transaction. When dependencies between groups are due to more than one transaction type, the dependencies are combined by adding the active components and measuring the union of passive components. Thresholds are then established judgementally to allow transactions to be categorized into four groups: (HH) meaning it is used frequently and accesses a large quantity of data, (HL), (LH) and (LL).

One constraint placed in Baker's model is that a logical distribution may not contain any HH dependencies. The dependency between two groups is then described using an orientation and three consolidated dependencies (HL, LH and LL). Values of these thresholds determine the number of groups in the logical distribution, the extent of the group's autonomy and the strength of the dependencies. Baker continues to discuss two types of dependency support, namely data communications when data currency and integrity are important or data duplication which yields good response time, management control and system availability.

Although Baker's approach may take several iterations to reach a satisfactory distribution, it highlights the importance of placing the data processing functions and associated data close to their users. Only in this way can a successful degree of nodal autonomy or an acceptably low level of interdependency be obtained, allowing the benefits of distributed processing to be utilized.

2.5.2 Theoretical Models

In contrast to Baker's trial and error approach, some research proposes theoretical methods to obtain logical distributions. Similar to the data allocation models, these approaches generally require heuristic techniques to make their use practical [NAVA84]. Many of these models are concerned with affinity among attributes and attribute clustering [CERI85],[NAVA84]. The work done by Navathe, Ceri, Wiederhold and Dou will be reviewed in some detail, so that an understanding of the clustering techniques can be achieved.

Navathe et.al. propose a set of algorithms to deal with the vertical partitioning problem (VPP).

The approach consists of two phases where in the first phase the design is independent of specific cost information. The second phase performs cost optimization from knowledge of a specific application environment. The model also deals with three environments for the vertical partitioning problem; a single site with one memory level, a single site with multiple memory levels and multiple sites [NAVA84]. The multiple site allocation is most relevant to this review.

The inputs to the model are the logical accesses of the transactions to the attributes (i.e., number of accesses to object instances for one occurrence of a transaction at a site) and the relevant design parameters such as cost of storage, access and transmission. The steps in the first phase are summarized below [NAVA84]:

1. *Construct Attribute Affinity Matrix (AA Matrix):* The objective of this function is to construct an AA matrix which records the affinity or imaginary bond between attributes. The affinity measure is based on the logical access information which has been obtained. This includes whether or not a transaction uses a particular attribute, whether the transaction is retrieval or update and the number of accesses to the object for one occurrence of the transaction. The affinity measure recorded is the sum of these accesses per time period (i.e., per day).
2. *Cluster the Attributes:* The objective of this function is to group the AA matrix so that attributes with high affinity are clustered together as are attributes with low affinity. This is accomplished through a heuristic algorithm that diagonalizes the AA matrix to produce blocks of jointly accessed data items.
3. *Partitioning:* The authors provide two mechanisms for partitioning. The first provides partitioning for non-overlapping fragments. This attempts to find the ideal location on the matrix to form the partition of two non-overlapping fragments, so that the fragments are balanced with respect to transaction load. In general terms, if 'n' is the number of attributes, 'n-1' points on the diagonal are considered. At each point, the matrix is "split" into an Upper and Lower fragment. A count is made of the total number of access of transactions that need only fragments in the current upper block (labeled (CU)) and current lower block (labeled (CL)), and a count of the number of transactions that need both fragments (CI). The goal then is to select the point in the matrix such that the goal function (z) is maximized: $\max z = CLxCU - CI^2$

The second mechanism allows for partitioning with overlapping fragments. This requires the use of two points on the diagonal, x1 and x2 where attributes between x1 and x2 constitute the intersection. The goal function is the same for the non-overlapping case however consideration is given to read only versus update requests. In the case of read only there are advantages to sharing the data among the two fragments whereas updates need to be directed to both fragments for consistency.

The process as presented is iterative, as it is most likely that the vertical partitioning will result in several fragments. In order to reduce the computational complexity of the problem, the authors chose a suboptimal approach whereby each application of the VPP produces two fragments. These fragments then become independent subproblems or fragments on which the algorithms are repeated to further split the fragments. This is repeated until no further benefits are gained [NAVA84].

The final phase of the methodology deals with the allocation of these fragments to sites. Four cost factors are considered in this stage. These include the cost of irrelevant attributes accessed within a fragment, cost of accessing fragments for retrieval and update, storage cost and transmission cost. These cost factors are then assigned weights according to their importance in the overall optimization model. In a distributed environment the transmission costs receive the highest weight in their model. To summarize, a table of partitions versus allocations is maintained and for each possible partition the algorithm attempts all possible fragment allocations (m^2 cases for m sites). The "least cost" pair is selected. In the case of replicated sites, a final algorithm is invoked. This algorithm looks at each fragment independently and allocates additional copies to sites until no further benefit is gained [NAVA84].

Other optimization models exist which address horizontal and replication issues [CER183a,CER183b]. These models analyze the logical distribution in terms of objects and links and the relationships between them. They assume explicit knowledge by the user in terms of potential uses of the database, transaction frequencies and cardinalities (number of instances) of objects and links, etc. Given these inputs they produce formal solutions to the design problem.

2.5.3 Summary

The methodologies and solution methods which have just been reviewed address the distribution problem from a higher level than the pure file allocation models. This is done by assessing the

dependencies between transactions or the affinity or bonding between attributes. These approaches highlight the many issues that must be taken into account during logical design. It is important to note that there does not appear to be any one model that addresses horizontal and vertical partitioning under one methodology, due to the complexity involved.

2.6 Related Work

As a final note, Hevner has pointed out the importance of combined research in the areas of data allocation and query optimization. As he states, the data allocation problem generally assumes a given data access pattern and the query optimization problem assumes a fixed data allocation. If these assumptions were generalized, designers could develop models that combine the allocation of files and the support of query processing in the most efficient manner [HEVN84].

Hevner reviews some of the recent work done in this area. The work of Elam minimizes the amount of data sent for a specified set of query processing strategies, under the constraint that one file is stored separately from others in the query to promote parallel processing. Another approach by Apers develops an algorithm to minimize data transmission costs by clustering relations. The constraint placed on this model is that only one copy of each relation is stored in the system (i.e., no redundancy) [HEVN84]. In agreement with Hevner, further research in this area is needed.

2.7 Conclusions

This chapter enumerated the many issues associated with the distribution of data and file allocation in a distributed environment. Several attempts to deal with these problems have been addressed in the literature. Early models focused solely on the optimal placement of files, with extensions added later to address program and hardware allocation in parallel. The major weakness of these models was the assumption of a file being the unit of allocation with no incorporation of file partitioning or file interdependence. Other methodologies and solutions were

presented that address these issues. These methodologies focused on transaction dependencies and bonding between attributes when analyzing the placement of data.

In summary, in a distributed environment a database designer will be faced with many different problems, including how to partition files, where to place these partitions and what degree of redundancy should be incorporated. Although an optimal solution to all of these problems is not realistic at this time, various models exist for different classes of these decisions. Integration of a solution to some of these problems into an interactive tool, would greatly reduce the complexity involved. Such a tool would provide a building block on which the designer could utilize his judgement, to design a feasible distribution scheme.

CHAPTER 3

DESIGN

3.1 Objectives

The objective of this project is to design a tool that can be used to aid the database designer in the distribution of data. Total optimality and automation of all aspects of the distribution problem still requires a great deal of research. A tool therefore can only expect to provide guidance and recommend solutions, as ultimately the design process still relies to a large extent on the designer. In that regard, the distribution algorithms developed should be integrated into an interactive design tool. With such a tool, the designer can review the steps that have been taken and modify the results as desired.

3.1.1 Design Overview

The tool being designed is an extension to two automated tools in existence today. The first is the "document_handler" program which processes application documents and determines the functional dependencies (FD's) they represent. The FD's are used as input to the "bern2" program, which produces a database in third normal form (3NF). The data distribution tool (hereafter referred to as "dist_data") then deals with the question, what if I now want to distribute the data ? A high level overview of the functionality of these tools is required in order to understand the interdependencies between them.

3.1.1.1 Document_handler Program

In distributed processing, one of the fundamental keywords is *document*. A collection of the documents used in an organization will tell you a great deal about the database schema required. A complete set of user's documents should supply all of the data items used by an organization. Problems arise though, as different user documents refer to the same data item by different names (synonyms) or the same name is applied to different data items on different documents

(homonyms). Another problem arises when analyzing these documents, as some data items names appear on a user document but are not actually stored in the database. These are derived data items that may be computed from other data values. The database designer must examine these documents to generate a nonambiguous list of data items for the database [WOEL71]. This can be an arduous task.

The document handler tool was developed to provide automated assistance in interpreting these documents. The tool takes as input all user documents and their associated columns, along with a specification of their use as input, resident or output documents. The documents and columns are scanned one by one, and the following activities take place [COHE79]:

- Deletion of synonym/homonym names
- Removal of insignificant columns (some columns may not be included in the database due to their nature. A signature field for example)
- Solve undeclared output

These activities are accomplished by interactive dialogue with the user. The system contains a sophisticated mapping between all documents and columns. This allows a cross reference listing of any column to any document as well as a listing of any document, its type and all columns it contains. This information is stored and maintained by the program.

Once a unified list of document names and attributes has been derived, keys are specified for each document. These keys are then used in conjunction with the columns in a document to form functional dependencies. These functional dependencies are then used as input into the bern2 program.

3.1.1.2 Bern2 Program

The bern2 program is a tool which automates the steps of Bernstein's algorithm, that produces a database schema which is in third normal form (3NF) with a minimal number of relations. A database in 3NF contains no extraneous attributes, no partial dependencies and no transitive

dependencies (a transitive dependency exists when you can get from a key to a non-key to a non-prime). The removal of extraneous attributes and redundant FD's is particularly applicable to the FD's introduced by the document_handler program, as many documents will use the same keys and contain a great deal of data items in common.

The steps of Bernstein's algorithm can be summarized as follows [COHE79]:

1. Eliminate from the functional dependencies those data items that can be derived from other functional dependencies (extraneous attributes)
2. Eliminate from the set those functional dependencies that can be derived from the remaining set of FD's (redundant FD's)
3. Group the remaining functional dependencies into sets with equivalent left hand sides
4. Merge the groups that have equivalent left hand sides (i.e., keys)
5. Remove transitive dependencies from the data items
6. Construct relations based on the groups of functional dependencies.

3.1.1.3 Data Distribution Tool

The dist_data tool is the focus of this project. It extends the work of the tools just reviewed, by addressing the fact that the resultant database may be used in a distributed environment. It utilizes much of the information obtained in the document_handler and bern2 programs. Specifically, the document_handler provides critical information about an application which is needed to make intelligent distribution decisions. This information includes application transactions in the form of input and output documents.

The bern2 program removes all extraneous attributes and transitive dependencies to produce a database in 3NF. These relations coupled with the information on customer documents provide the basis for defining meaningful partitions and file placement onto nodes in a network.

As previously discussed, there are two main concerns when dealing with data distribution. The first is to evaluate the utility of partitioning data objects into fragments, and the second is, once partitioned how these fragments are allocated to nodes on a network. This tool addresses both these issues. The intent is not to promise true optimality, but to provide a feasible solution to

these complex issues.

3.1.1.4 Limitations

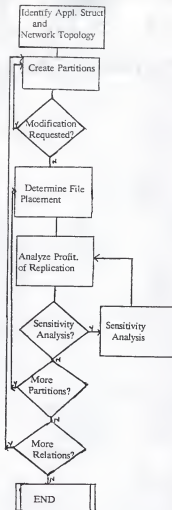
Some of the requirements originally outlined for an automated document handler [WOEL81], were not implemented in prototype document_handler program. Specifically, an indication of the frequency of use on a per document basis and the ownership of data items (i.e., what document(s) own each particular data item). These pieces of information are required in the distribution algorithms presented. As these modifications would not prove difficult to make, they will be assumed as input for the purpose of this project.

3.2 Detailed Design

This section presents the detailed design of the dist_data tool. The design algorithm involves four basic steps that can be summarized as follows:

1. Identify application structure and network topology
2. Create Partitions/Files
3. Determine File Placement
4. Analyze Profitability of Replication (Branch and Bound)
5. Sensitivity Analysis

The steps of the algorithm are iterative in nature; step 2 is repeated for all relations produced from bern2, and steps 3, 4 and 5 are repeated for each file or partition generated. Figure 3.1 provides a flowchart of this activity. The remainder of this chapter will discuss each step of the design algorithm in detail.



Flowchart of Dist_data Process
Figure 3.1

3.2.1 Step 1 : Identify Application Structure and Network Topology

Data distribution cannot begin without understanding the organizational structure and distribution scheme for any particular application. The first step in the design process is to gather this information from the user. For ease of use the tool will prompt the user for the needed information, which includes:

1. Number of nodes on the network
2. Identification of organizational groupings (i.e., INVENTORY, BILLING, etc.)
3. On a per organization basis:
 - Which node the organization resides on
 - A list of nodes which interconnect with this organization and the associated communication costs of each interconnection
 - A list of the documents used by the organization

With this information in hand, several required pieces of information can be compiled. First, the application network topology has been defined. The system may now build cost tables which reflect the transmission costs incurred for queries and updates. It is important to point out one requirement of the system which is needed by later algorithms. The file placement algorithm (step 3) and replication algorithm (step 4) require cost figures to be associated between any two pairs of nodes on the network (i.e., a fully connected network). The problem can be demonstrated as follows. Consider a system with three nodes, where nodes 1 and 3 are connected to node 2, but no connection exists between nodes 1 and 3. This topology is illustrated in figure 3.2.



Example Network Topology
Figure 3.2

The cost tables which need to be built require that a cost be associated with nodes 1 and 3. In the example above, this can be accomplished by taking the sum of the communication costs between nodes 1 and 2, and nodes 2 and 3. In a situation where multiple paths are available, the least cost path should be associated with the "imaginary" link. The system should detect these missing links and compute the cost factor. Once a cost is associated with all of the nodes in the network, the cost table can be built. Figure 3.3 shows an example of a cost table for three nodes. Reading the table across tells you that the cost of communication between node 1 and itself is zero, node 1 and 2 is 8 and between nodes 1 and 3 is 12. This information is required in steps 3, 4 and 5 of the algorithm.

(Nodes)

	1	2	3
1	0	8	12
2	8	0	20
3	12	20	0

Network Cost Table

Figure 3.3

The other piece of information obtained from the user deals with the application structure. The system now knows on which nodes each document or transaction resides. This node identifier must be stored in conjunction with the document entries as it is vital in determining a meaningful partitioning of relations (this will become clear after reviewing the algorithm for step 2).

The information obtained in this step serves as a framework for all other steps. It should be noted that this is the only step which *requires* information from the user (all other steps of the algorithm make use of the information provided by the document_handler and bern2 programs). Appropriate error detection and recovery will be provided when parsing the user input. This includes syntactic checks as well as checks for invalid document names or invalid interconnect nodes.

3.2.2 Step 2: Create Partitions/Files

Step 2 of the algorithm is a key step in the distributed design process. In this step, the relations produced by the bern2 algorithm are analyzed to determine if partitioning should be applied. Recall that initially the application has a list of documents which they own and execute. The document_handler program removes all homonyms, synonyms and insignificant columns to insure that the data in the document is interpreted correctly. The bern2 program then takes these documents (represented as FD's) and produces a 3NF database. These relations do not reflect customer usage however, as one relation may now contain pieces of data from several different documents. Chapter two reviewed the importance of analyzing the application's use of data , particularly in regard to the dependency which exists between files. These dependencies may reveal clusters of data which belong together or that could form meaningful partitions. Once these partitions are identified, the file allocation algorithms can deal with the physical placement of these fragments onto sites in a network.

Several different approaches which have been taken in dealing with partitioning were discussed in Chapter two. The heuristic chosen for this design process, is to analyze the frequency of use of each attribute in the relation on a per node basis, and cluster these attributes together into a file. This approach, although simplified is very similar to the approach used by Baker. The justification for this algorithm is as follows:

- The grouping of attributes by frequency of use agrees with the underlying message in the literature; store the data where it is most frequently used.
- Partitioning the data in this manner usually results in partitions being stored at the same node as the owner of those attributes. This proves beneficial, as most of the time people are "greedy" in the sense that they want to maintain control over the data they own.
- Availability of information: The document_handler program contains all the needed information to analyze the relations in this manner. Each document entry contains all related data items and the frequency of use. Step 1 of the distribution algorithm also associated a node with each document. All of this information is readily available and feasible for the scope of this project.

In order to accomplish the partitioning, three steps are needed. These steps build a document table associated with the relation and sort it by node, analyze the data frequencies and then form partitions based on this frequency. At the conclusion of this step the designer may review the partitioning which has taken place, and make modifications if desired. These functions work on one relation from the bem2 algorithm at a time. They are therefore repeated until all relations have been processed.

The phases of this step are reviewed in detail below.

3.2.2.1 Build Document Table

Figure 3.4 and 3.5 contain one relation produced by bem2 and three document entries as they would look after step 1 of the distribution process.

REL1 (KEY_ATTR, ATTR1, ATTR3, ATTR5)

Bem2 Relation

Figure 3.4

<u>NODE</u>	<u>FORM</u>	<u>KEY_ATTR</u>	<u>ATTR1</u>	<u>ATTR2</u>	<u>ATTR3</u>	<u>ATTR4</u>
1	form1	10	10*	10	10	10*

<u>NODE</u>	<u>FORM</u>	<u>KEY_ATTR</u>	<u>ATTR5</u>	<u>ATTR8</u>
3	form2	10	10	10

<u>NODE</u>	<u>FORM</u>	<u>KEY_ATTR</u>	<u>ATTR3</u>	<u>ATTR8</u>
3	form3	5	5	5

Document Entries

Figure 3.5

The first step is to build a document table which corresponds to the relation being examined. In this example, only those documents using the key attribute and attributes 1,3 and 5 are of

concern. This table is then sorted by node for ease of processing in the next step. Figure 3.6 shows the document table at the conclusion of this step. Note that attributes 2,4 and 8 used by the forms, do not appear in this table.

NODE	DOC	KEY_ATTR	ATTR1	ATTR3	ATTR5	ATTR7
1	form1	10	10*	10		
3	form2	10			10	
3	form3	5		5		

Document Table
Figure 3.6

3.2.2.2 Analyze Frequency of Use

In this step, each data item or attribute is examined individually and an aggregate usage at each node is determined. Two important assumptions are included in this step:

1. Keys are not analyzed, as the key to this relation is required in any partition created (i.e., the key must be duplicated for access)
2. An access to an attribute by the owner is taken as an update transaction (recall that an owner(s) of a data item must be indicated. This is represented by a '*' in these examples. It is possible for more than one owner to exist for any given attribute.) These accesses represent exception rates such as new customers being entered, etc.

Using the example in Figure 3.6, the following data would be compiled:

- ATTR1: 10 updates node 1
- ATTR3: 10 queries node 1, 5 queries node 3
- ATTR5: 10 queries node 3

3.2.2.3 Create Partitions

This step processes the frequencies above to determine the partitioning. In most cases the aggregate total of query and update requests from a node is used to determine the owner. In the case where a query count from one node equals the update count from another node, the update access is given the higher priority (It is recognized that further heuristics could be applied when

queries are very close to updates, etc. Further heuristics are not addressed in this design). Each attribute is therefore associated with the node which contains the highest frequency of use. The final step combines all attributes from the same node into one file.

The final partitions resulting from the example above are:

```
FILE1 = (KEY_ATTR, ATTR1, ATTR3)
FILE2 = (KEY_ATTR, ATTR5)
```

3.2.2.4 Query User

After the partitions have been formed, the results should be displayed to the user. At this point the user will be allowed to modify the partitioning if desired. Three operations will be allowed:

1. MOVE ATTR ____ to FILE ____
2. CREATE FILE ____ (This will create a new partition, which may then be populated through a series of MOVE requests)
3. MERGE FILE ____ and FILE ____

These operations will again be prompted for, so that minimal effort is required on the user's behalf.

3.2.3 Step 3: Determine File Placement

Once the relation has been partitioned into files, each of these files needs to be analyzed individually to determine where on the network they should be placed. The objective function chosen is to place the file at the location which minimizes the overall communication costs in regard to update and query requests. To accomplish this, transaction tables must be built to indicate the frequency of requests issued against this file. Query and update tables are then built and analyzed to determine the optimal placement of the file.

3.2.3.1 Build Transaction Table

A transaction table needs to be built for each partition or file produced in step 2. The table contains a count of the number of queries and updates issued against this file from each node.

This information was previously gathered for all attributes. This step however is only concerned with those attributes associated with the partition being analyzed. The same rules for distinguishing between query and update requests in step 2, also apply to this step. In addition, one other assumption is made:

- It is assumed that all attributes in a document/file are accessed together (i.e., cannot access individual attributes out of a file, the entire file is retrieved). This implies that if a form has some queries and some updates against the relation, an update transaction is assumed. It would not seem appropriate to consider these as two separate accesses, so one access is assumed with updates weighing more heavily.

The output of this process will be a transaction table which holds the associated access frequencies.

Figure 3.7 shows an example of a transaction table (an additional node has been added from previous examples).

Nodes	Query	Update
1	10	10
2	30	10
3	15	0

Transaction Table
Figure 3.7

3.2.3.2 Build Cost Table

Cost tables must now be created to reflect the query and update costs that would result from placing the file on any given node. These costs are computed using the transaction table and the network cost obtained in step 1.

For example, assume the following entries exist in the transaction table and network cost table for Node 1:

Transaction Table			Network Cost			
Node	Query	Update		1	2	3
1	10	10	1	0	8	12
2	30	10	:			
3	15	0	:			

The query costs associated with placing the file at node 1 can be computed as follows:

1. Node 1 to itself: cost = 0
2. Node 2 to node 1: cost = 240 (30 query requests from node 2 at a cost factor of 8)
3. Node 3 to node 1: cost = 180 (15 query requests from node 3 at a cost factor of 12)
4. Total query costs incurred = 420

The same algorithm is used to determine update costs. The total cost associated with each node is the sum of query and update costs. These costs are computed for every node on the network. The node representing the minimum cost is chosen for file placement.

3.2.4 Step 4: Determine Profitability of Replication

The next step of the distribution is to determine the profitability of replication. This is accomplished by building a decision tree and performing a branch and bound search.

In theory, the branch and bound search looks for an optimal solution by defining initial upper and lower values of the objective function (in this case the objective function is the minimization of communication costs). From the feasible solutions, the best solution is made the upper value (U) of the problem. All other solutions are matched against this solution in an attempt to find a better solution. Any solution which produces a value higher than (U) are deleted, as further branching would not lead to a better solution. This process continues through a series of iterations in an attempt to find the optimal solution [GREEN78].

The branch and bound search technique was reviewed briefly in Chapter two. This discussion highlighted the fact that even though this search procedure bounds the number of required computations, heuristics are still needed to reduce the potential computational requirements. In this light, the heuristic identified for this design is a simple "greedy" heuristic. At any given point in the search, only the best solution is kept. Although it is realized that this can not guarantee the true optimal placement of the file, it provides a reasonable placement within the scope of this project.

To correlate the above discussion to the file allocation problem, consider the placement of a file on a network consisting of three nodes. Three levels of the decision tree are depicted in figure 3.8 for purposes of this discussion.



Partial Decision Tree
Figure 3.8

At the first level, each vertex represents a file assignment to a given node (denoted by 1's in those positions corresponding to file nodes, 0's elsewhere). The file placement algorithm discussed in step 4, determines the optimal placement of one copy of the file to the network. For this example, assume this is node 1 (represented by a 100 at level 1 on the graph). With the "greedy"

heuristic applied, the search will now only be concerned with replication schemes that involve node 1. The search proceeds to level 2, where the replication choices are a copy of the file at nodes 1 and 2, or a copy at nodes 1 and 3. The cost associated with each of these choices is analyzed. Obtaining the cost figure is very similar to the file placement algorithm; however, it is modified slightly to account for the replication. For example, when analyzing the costs associated with file copies allocated to nodes 1 and 2 the following changes will be made. The cost associated for queries at both node 1 and node 2 are zero, as each node will now have a local copy. A query from node 3 can now be satisfied from either node 1 or 2, so the lower cost is assumed. Update costs are increased however, as updates must be propagated to all file copies. In this case, the update cost associated with the replication would be the sum of update costs associated with nodes 1 and 2. Again, the total cost is the sum of query and update costs. If this cost is less than the cost associated with only 1 copy of the file, then this choice is now considered the "best" placement and the search continues to level 3 (which represents total replication in this example). If the costs associated with replication is higher, the search has ended (after both nodes are tested).

3.2.5 Step 5: Sensitivity Analysis

As the replication decision may be far from optimal, a sensitivity analysis is being provided. The intent of this step is to allow the designer to modify the communication costs associated with a given topology and re-analyze the file placement. In other words, the designer will be able to request a reiteration of the branch and bound search with the communication cost increased by some specified amount. If the result after this run differs significantly from the first, this may indicate that replication is desirable in either case. For example, if a 5% increase in communication costs shows that replication is desirable this would indicate that the file placement is very sensitive. This type of sensitivity may imply that replication is in order.

At the conclusion of this step, the file allocation process has been completed for one file or partition. Each original relation may be split into several partitions. Steps 3, 4 and 5 (if desired) will need to be repeated for each partition.

3.3 Summary

This chapter presented the design for the `dist_data` tool. This tool takes the information regarding customer documents coupled the 3NF relations they represent, and automates a distribution analysis. This analysis includes recommended partitioning for the relations based on file dependencies and frequency of use, and the allocation of these partitions onto the nodes in a network. The allocation scheme is based on minimizing communications costs for update and query transactions. The process is iterative, examining one relation at a time and then each resulting partition in that relation. This process is repeated until all relations have been analyzed.

The interface to the tool is interactive in nature, allowing the designer/user to make modifications if desired to appropriately steer the distribution design. A minimum amount of error detection and recovery has been provided.

Chapter four contains a detailed example of the `dist_data` operation.

CHAPTER 4

EXAMPLE DESIGN

4.1 Example Design/Results

This chapter presents a complete design, following the steps outlined in Chapter three. The system used in this example, models a manufacturing firm which consists of three organizations: CUST_REL, ORDER/BILLING and INVENTORY. These organizations will be distributed across a three node network.

The example starts by showing a logical listing of the documents, as they would appear at the end of part one of the document_handler program. Appendix I contains a list of these documents. Each document is labeled with a type (i.e., INPUT, RESIDENT or OUTPUT) and has a document key specified. An owner is also indicated for each data element (denoted by a '*').

Appendix II contains the FD's as produced by the document_handler. These FD's are used as input into the bern2 program. At the conclusion of the bern2 run, two 3NF relations are produced.

The dist_data tool commences after the results from the above programs have been obtained. Appendix III contains a step-by-step example of the procedures used. The remainder of this chapter will summarize the processing that occurs at each step.

The first step of the dist_data tool, is to identify the application structure and network topology ((III.1). This notation will be used throughout the chapter for Appendix III, point 1). Information is solicited from the user regarding the organizational structure, the documents used and the network topology. Note that in the example, the system must detect the missing node interconnection between nodes 2 and 3, and search for the least cost connection (in this case there is only 1 path between nodes 2 and 3). At the conclusion of this step, a network cost table has

been built and each node has a list of associated documents that are resident at the node.

The next operation examines the utility of partitioning (III.2). This step of the algorithm deals with one relation at a time, which in this case is REL1 produced by bern2. The table is constructed by identifying all documents which use the attributes associated with this particular relation. The owner of a data item is indicated by a "*", and the numbers in the table represent the frequency of use. These frequencies are then analyzed on a per data item basis to determine the owner (i.e., node) of that data item (denoted by an 'X'). All data items with a common owner are then merged to form a partition. In this example, Node 1 exhibited the highest frequency of use for all data items and therefore no partitioning was done. At this point the designer may display the partitioning that has been created and perform modifications if desired (see section 3.2.2.4).

Step 3 of the algorithm determines the placement of all partitions/files produced. It deals with one partition at a time, so it is repeated for each partition produced in step 2 (in the example so far, only 1 partition has been created). Step 3 begins by building a transaction table (III.3). This table represents the aggregate usage of this file from all nodes. Recall from section 3.2.3.1, there are two critical assumptions used when building this table. The first is that access by an owner is interpreted as an update transaction. Secondly, a document or file is assumed to be accessed in entirety. The importance of these assumptions can be highlighted by examining the HIST_FILE document. This document accesses c_no and part_no as a query request and tot_price as an update. This transaction from node 1 is interpreted as an update request of frequency 30. The remainder of the requests from node 1 are query only, so the table is populated with 30 query requests and 30 updates.

Once the transaction table has been built, the cost for query and update resulting from placing the file on any given node, can be determined. The network cost table (from step 1) is used in

combination with the transaction table, to compute these costs. For example, Node 1 issues 30 queries and 30 updates. If this file were placed on Node 1, there would be no cost for queries or updates from Node 1, as all the information is local. However, there are query and update requests to this file from Nodes 2 and 3. Node 2 issues 6 queries and 5 updates with a cost factor of 8. Reading the query, update and network tables down the Node 1 column, this means if the file were placed at Node 1, a cost of 48 would be incurred in queries and 40 in updates from Node 2. The computation continues in this manner in order to analyze the total cost from any node. The node which represents the minimum cost is recommended for the file placement (Node 1, in this example).

After a single file placement is recommended, the tool looks at the profitability of replication (III.4). This represents step 4 of the distribution algorithm. As previously discussed in section 3.2.4, the tool progresses only with the "best choice" at any given time. From the preceding step, node 1 is chosen for the file placement. Replication options at level two of the cost graph include two copies of the file, at either nodes 1 and 2, or node 1 and 3. The costs are recalculated under the assumption of replication, to determine if either option results in a lower cost. In the example, replication proves profitable at nodes 1 and 3. Therefore, the final recommendation for file 1 is to replicate the file at nodes 1 and 3, at a total cost of 572.

The final step of the algorithm (step 5) allows the designer to perform sensitivity analysis (III.5). This has not been illustrated in the example, but would involve receiving a cost increase parameter from the user (i.e., 5%), and recomputing the cost tables and reanalyzing the distribution.

Appendix III.6 displays the network after the placement of file 1. The communication traffic resulting from this placement is also illustrated. In analyzing this placement, there are two points of interest. A general view taken by [ROTH81] discussed in Chapter two, stated that in reality queries tend to be large and complex whereas updates tend to be small and simple. A file

placement that minimizes query traffic is therefore beneficial. The file placement produced in this algorithm has modeled this assumption. The major query requests have been fulfilled at both nodes 1 and 3. Updates do incur communication expense but these updates are small (part_no and tot_price). The frequency of access to this file at Node 2 is so low, that further replication would not prove advantageous. The second theory discussed in the literature, is that a single copy of a file suffices under most conditions, if the ratio of query to updates is more than or equal to 50% [CASE72]. Looking back at the transaction table for file 1, the ratio of query to update (77 queries, 37 updates) is 48%. The results of this algorithm have stayed within these guidelines (it should be noted, that this comparison holds throughout the remainder of this example. All file placements meet these guidelines).

Processing now returns to step 2 of the algorithm, to begin analysis of the next bern2 relation (III.7) At the conclusion of this step, REL2 has been partitioned into two files.

Step 3 of the algorithm: determine file placement, is once again invoked for File 2 (recall that this step analyzes a single partition at a time). The processing steps are identical to that of file 1; the node that minimizes the total costs is found (III.8) and then replication is analyzed (III.9). The final recommended placement for file 2, is again two copies of the file residing at nodes 1 and 3. The total cost for this placement is 264. III.11 shows the overall view of the network with both files 1 and 2 allocated.

To complete the processing, one partition is left to analyze (File 3). The algorithm returns to step 3, to determine the single file placement for this partitioning (III.12). Again, the processing is identical to that of file 1 and 2 and is illustrated in III.12 - III.14. In the case of file 3, replication is not profitable and therefore a single copy of the file is placed at node 3. The cost associated with this placement is 0.

Appendix III.15 shows the final results after the dist_data program has processed all relations and partitions. To summarize, file 1 and 2 are replicated on nodes 1 and 3, and a single copy of file 3 exists at node 3. The overall cost resulting from these allocations is 836.

CHAPTER 5

CONCLUSIONS

5.1 Discussion of Results

The goal of the system presented in this paper is to *assist* the database designer in finding solutions to the problems of file partitioning and file placement in a distributed system. The basic premise of the system is that user documents contain the needed information to make reasonable distribution decisions. A complete set of user documents supply the data items used by an organization and model the transactions (in the form of input and output documents). This information, coupled with an understanding of the ownership of these documents, allows file partitioning and file placement to be made based on the frequency of use and ownership of the data items. By automating the analysis of this data and proposing solutions, a great deal of complexity is removed from the designer. This is not to suggest that this is the optimal solution nor the only solution to the problem, as there are many issues which can impact these design choices. As such, the system was designed to be interactive in nature, allowing the designer to make modifications as required to appropriately steer the design process.

Research continues to address the key problems of file partitioning, file placement and redundancy considerations, in an effort to find optimal solutions. The underlying message in all these solutions is clear however; a major determinant in the increased use of distributed processing will be an increased ability to get data where it is needed, and used most frequently.

5.2 Extensions

Numerous enhancements could be made to increase the benefits derived from this system. These include extensions to the partitioning and file allocation heuristics, as well as the overall user interface to the system.

The partitioning scheme developed for this system does not directly address partitioning by occurrence, or horizontal partitioning. Although the system will handle the same organizations existing at multiple nodes and maintaining ownership of the same documents, it will assign the file to one of the two nodes only. For example, there is nothing to preclude the user from entering the ORDER/BILLING organization on two different nodes, each owning the same documents. If the frequency of use were the same against each document, the system will arbitrarily designate one node as owner. The single file allocation of this file would result in the node that represents the minimum communication cost. In this case the algorithm will most often conclude that replication is profitable at the second node, unless the communication costs vary significantly. The algorithm will not however, correctly address the other remote requests to this file, as it assumes both files are identical. The second problem arises when the frequency of access varies between the two locations (i.e., more activity with certain part numbers than others). In this case, the algorithm will cater to the node which exhibits the higher access rate. These problems point out the importance of allowing input from the designer, as the designer can easily remedy this situation by creating and allocating new partitions. A more sophisticated method of detecting and dealing with horizontal partitioning would be desirable however.

Another specific enhancement that could be applied to the partitioning algorithm, is to supply an improved heuristic for differentiating between updates and queries. The current design gives weight to an update in the case where query and update rates are identical, but does nothing more. It would be desirable to further analyze this difference and extend the heuristic (i.e., if there were only 10% more queries than updates, should the algorithm still favor updates?).

Other limitations of the system exist in the file allocation and replication heuristics. The model is very restrictive, in that it only considers query and update communication costs. No consideration is given to such things as storage restrictions and costs, communication channel load, etc. Chapter

two highlighted several different models which take various system resources into consideration. The integration of any one (or more) of these issues into the model would certainly provide a more refined file allocation.

Finally, limited attention was given to the design of the user interface to the tool. A primitive method of displaying and modifying the information was described, that although is sufficient, could be greatly enhanced. One could envision a menu driven system with graphical display capabilities, in which the designer could control the total execution of the system, requesting any of the steps to be run in any order. The graphics would allow a cleaner display of the processing at any given step, or a visual view of the network topology and file placement.

- [BAKE] Baker, Charles T., "Application Analysis for Nodal Autonomy", International Business Machines Corporation, Poughkeepsie, New York.
- [BRAY81] Bray, Olen H. "Distributed Database Design Considerations", *IEEE Tutorial on Distributed Processing*, Third Edition, IEEE Computer Society Press, 1981 (ppgs 451-464)
- [CASE72] Casey, R.G., "Allocation of Copies of a File in An Information Network", *AFIPS Conference Proceedings*, Vol. 40, Spring 1972 (ppgs 617-624)
- [CER183a] Ceri, Stefano, Navathe, Shamkant and Wiederhold, Gio "Distribution Design of Logical Database Schemas", *IEEE Transactions on Software Engineering*, Vol. SE-9, No. 4, 1983 (ppgs 487-503)
- [CER183b] Ceri, S. and Navathe, S.B., "A Methodology for the Distribution Design of Databases", *Compccon*, Spring 1983 (ppgs 426-431)
- [CER185] Ceri, S. and Pernice, B., "DATAID-D: Methodology for Distributed Database Design", *Computer-Aided Database Design: The DATAID Project*, A. Albano, V. DeAntonellis, and A. Pileva (editors), Elsevier Science Publishers B.V., 1985
- [CHAM81] Champine, G.A., "Six Approaches to Distributed Databases", *IEEE Tutorial on Distributed Processing*, Third Edition, IEEE Computer Society Press, 1981 (ppgs 480-483)
- [CHU69] Chu, Wesley W., "Optimal File Allocation in a Multiple Computer System", *IEEE Transactions on Computers*, Vol C-18, No. 10, Oct. 1969 (ppgs 885-889)
- [COHE81] Cohen, Meir, "A System for Automatic Generation of Relational Databases", *A Master's Report*, Kansas State University, 1981
- [DRAF80] Draffan, I.W. and Poole, F., "The Classification of Distributed Database Management Systems", *Distributed Databases*, Cambridge University Press, 1980 (pgs 57-81)
- [GAVI86] Gavish, Bezalel and Pirkul, Hasan, "Computer and Database Location in Distributed Computer Systems", *IEEE Transactions on Computers*, Vol. c-35, No. 7, 1986 (ppgs 583-589)
- [GROS80] Gross, J.M., Jackson, P.E., Joyce, J and McGuire, F.A., "Distributed Database Design and Administration", *Distributed Databases*, Cambridge University Press, 1980 (ppgs 285-322)
- [HEVN84] Hevner, Alan R., "Data Allocation and Retrieval In Distributed Systems", *Advances in Data Base Management*, Vol. 2, E.A. Unger, P.S. Fisher and J. Slonim (editors), Wiley Heyden Ltd, 1984 (ppgs 225 -252)
- [IRAN79] Irani, Keki B; Khabbaz, Nicholas G, "A Model for a Combined Communication Network Design and File Allocation for Distributed Databases", *1st International Conference on Distributed Computing Systems*, Computer Society Press, 1979 (ppgs 15-21)
- [LEVI79] Levin, Dan K. and Morgan, Howard Lee, "Optimizing Distributed Databases - A Framework for Research", *Distributed System Design*, IEEE Computer Society Press, Oct. 1979 (ppgs 321-326)

- [MAHM76] Mahmoud, Samy and Riordon, J.S., "Optimal Allocation of Resources in Distributed Information Networks", *ACM Transactions on Database Systems*, Vol. 1, No.1, 1976 (ppgs 483-497)
- [MAR184] Mariani, M.P., "Distributed Data Processing: Technology and Critical Issues", *TRW Series on Software Technology*, Vol. 4, 1984 (217 ppgs)
- [MURO85] Muro, Shojiro; Ibaraki, Toshihide, Miyajima, Hiachiro and Hasegawa, Toshihara, "Evaluation of the File Redundancy in Distributed Systems", *IEEE Transactions on Software Engineering*, Vol. SE-11, No.2, 1985 (ppgs 199 -204)
- [NAVA84] Navathe, Shankant; Ceri, Stefano; Wiederhold, Gio and Dou, Jinglie, "Vertical Partitioning Algorithms for Database Design", *ACM Transactions on Database Systems*, Vol. 9, No. 4, 1984 (ppgs 680-710)
- [PEEB81] Peebles, Richard and Manning, Eric, "System Architecture for Distributed Data Management", *IEEE Tutorial on Distributed Processing*, Third Edition, IEEE Computer Society Press, 1981 (pgs 451-464)
- [PURK83] Purkayastha, S., Kar, G., Berelian, E., Wong, P., Cascy, R.L., Farmer, L., Lo, P. and Chen, D., "Designing a Database Management System for Distributed Real Time Engineering Applications", *IEEE Transactions on Computers*, 1983 (ppgs 432-439)
- [ROTH81] Rothnie, J.B. and Goodman, N., "A Survey of Research and Development in Distributed Database Management", *IEEE Tutorial on Distributed Processing*, Third Edition, IEEE Computer Society Press, 1981 (pgs 484-498).
- [SCHN84] Schriederjans, Mark J., *Linear Goal Programming*, Petrocelli Books, Inc., 1984 (pgs 67-109).
- [WOEL81] Woelk, Darrell W., "The Generation of Entity-Relationship Diagrams from User Documents", *A Master's Report*, Kansas State University, 1981.
- [YU83] Yu, C.T., Siu, M.K., Lam, K. and Chen, C.H., "Adaptive File Allocation in Star Computer Network", *IEEE COMPSAC*, Nov 7-11, Computer Society Press, 1983 (ppgs. 537-545)

APPENDIX I
EXAMPLE DESIGN: DOCUMENT_HANDLER OUTPUT

Each DOCUMENT is labeled with a type (INPUT, OUTPUT, RESIDENT), frequency of use (i.e., per day), and has a document key

Each DATA_ELEMENT has an owner(s) associated with it (indicated by '*').

NEW_CUST	Document Attributes: Location.output Frequency: 5
*cust_no (doc.key)	
*c_name	
*c_str	
*c_cty	
*c_sta	
*c_zip	
*c_ph	

CUST_FILE	Document Attributes: Location.resident Frequency: 5
c_no (doc.key)	
c_name	
c_str	
c_cty	
c_sta	
c_zip	
c_ph	

CUST_LIST

Document Attributes:
Location.output
Frequency: 1

c_no (doc.key)
c_name
c_str
c_cty
c_sta
c_zip
c_ph

NEW_PART

Document Attributes:
Location.input
Frequency: 2

*part_no (doc.key)
*part_desc
*part_price

PART_SUP

Document Attributes:
Location.input
Frequency: 10

part_no (doc.key)
part_desc
*qty_made
*made_date
part_price

CALC_INV

Document Attributes:
Location.resident
Frequency: 30

part_no (doc.key)
qty_made
made_date
*made_to_date
qty_ord
ord_date
*ord_to_date
*qty_rem

INV_REPT

Document Attributes:
Location.output
Frequency: 1

part_no (doc.key)
part_desc
made_to_date
ord_to_date
qty_rem

CUST_ORD

Document Attributes:
Location.input
Frequency: 20

c_no (doc.key)
c_name
c_str
c_cty
c_sta
c_zip
c_ph
part_no
*qty_ord
*cord_date

HIST_FILE

Document Attributes:
Location.resident
Frequency: 30

c_no (doc.key)
part_no (doc.key)
cord_date (doc.key)
qty_ord
part_price
*tot_price

INVOICE

Document Attributes:

Location.output

Frequency: 10

c_no (doc.key)

c_name

c_str

c_cty

c_sta

c_zip

c_ph

part_no

part_price

qty_ord

tot_price

ord_date

APPENDIX II

EXAMPLE DESIGN: FUNCTIONAL DEPENDENCIES/BERN2 OUTPUT

** FUNCTIONAL DEPENDENCIES CREATED FROM DOCUMENT_HANDLER **

cust_no --> c_name, c_str, c_cty, c_sta, c_zip, c_ph

part_no --> part_desc

part_no --> part_desc, qty_made, made_date, part_price

part_no --> qty_made, made_date, made_to_date, qty_ord, ord_to_date, qty_rem

part_no --> part_desc, made_to_date, ord_to_date, qty_rem

c_no --> c_name, c_str, c_cty, c_sta, c_zip, c_ph, part_no, qty_ord, cord_date

c_no, part_no, cord_date --> qty_ord, part_price, tot_price

c_no --> c_name, c_str, c_cty, c_sta, c_zip, c_ph, part_no, part_price, qty_ord tot_price,
cord_date

** 3NF RELATIONS PRODUCED BY BERN2 **

REL1 : c_no --> c_name, c_str, c_cty, c_sta, c_zip, c_ph, part_no, tot_price

REL2 : part_no --> qty_made, made_date, qty_ord, cord_date, part_desc, made_to_date,
ord_to_date, qty_rem, part_price

APPENDIX III
EXAMPLE DESIGN: DIST_DATA

1) Identify Application Structure and Network Topology
(Algorithm Step 1)

PART A:

Enter No. Nodes on Network: 3

Enter No. of Organizations: 3

Enter Org. Name:

ORDER/BILLING on Node?

Enter Interconnect Nodes & Cost:

ORDER/BILLING

1

2,8

3,12

—8—

INVENTORY

3

1,12

—8—

—8—

CUST_REL

2

1,8

—8—

—8—

Enter Documents Owned:

CUST_ORD

HIST_FILE

INVOICE

NEW_PART

PART_SUP

CALC_INV

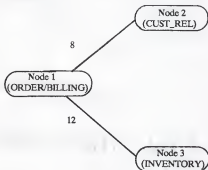
INV_REPT

NEW_CUST

CUST_FILE

CUST_LIST

Part B: Logical View of Network:



Part C: Build Cost Tables

(Nodes)

	1	2	3
1	0	8	12
2	8	0	20
3	12	20	0

2) Create Partitions/Files: Relation 1

(Algorithm Step 2)

Part A: Build Document Table (Sort by Node)

	c	c	c	c	c	c	c	part	tot
Node: Doc	no	name	str	cty	st	zip	ph	no	price
1 CUST_ORD	20	20	20	20	20	20	20	20	
1 HIST_FILE	30							30	30*
1 INVOICE	10	10	10	10	10	10	10	10	10
2 NEW_CUST	5*	5*	5*	5*	5*	5*	5*		
2 CUST_FILE	5	5	5	5	5	5	5		
2 CUST_LIST	1	1	1	1	1	1	1		
3 NEW_PART								2*	
3 PART_SUP								10	
3 CALC_INV								30	
3 INV_REPT								1	

Part B: Analyze Frequency of Use, Determine "Owners"

c_name

c_str

c_cty

c_st

c_zip

c_ph: 30 queries, 0 updates - Node 1 (X)

6 queries, 5 updates - Node 2

part_no: 60 queries, 0 updates - Node 1 (X)

41 queries, 2 updates - Node 3

tot_price: 10 queries, 20 updates - Node 1 (X)

Part C: Merge Attributes by Owner to form Partitions:

File 1 = (c_no, c_name, c_str, c_cty, c_st, c_zip, _ph, part_no, tot_price)

(No partitioning done for this file)

Part D: Modify Partitions (if desired)

3) Determine File Placement (Algorithm Step 3): Relation 1, File Name = 1

PART A: Build Transaction Table

	Queries	Updates
1	30	30
2	6	5
3	41	2

PART B: Build Cost Tables

(Network Cost Table from Previous Step)

(Nodes)

	1	2	3
1	0	8	12
2	8	0	20
3	12	20	0

Queries

	1	2	3
1	0	240	360
2	48	0	120
3	492	820	0
	540	1060	480

Updates

	1	2	3
1	0	240	360
2	40	0	100
3	24	40	0
	64	280	460

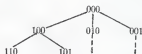
Total Costs: Node 1: $540 + 64 = 604$

Node 2: $1060 + 280 = 1340$

Node 3: $480 + 460 = 940$

***** Recommended File Placement at Node 1 ****

4) Determine Profitability of Replication (Algorithm Step 4): File 1



a) Node 1 and 2:

<u>Query Cost:</u>	<u>Update Cost:</u>
Node 1,2 = 0	Node 1 = 240
Node 3 = 492	Node 2 = 40
	Node 3 = 64
<hr/>	<hr/>
Total = 492	Total = 344

Total Cost = 492 + 344 = 836 (No replication)

b) Node 1 and 3:

<u>Query Cost:</u>	<u>Update Cost:</u>
Node 1,3 = 0	Node 1 = 360
Node 2 = 48	Node 2 = 140
	Node 3 = 24
<hr/>	<hr/>
Total = 48	Total = 524

Total Cost = 48 + 524 = 572 ** Replicate **

c) Continue to Branch, Nodes 1,2,3:

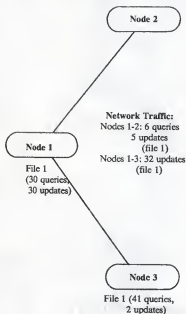
<u>Query Cost:</u>	<u>Update Cost:</u>
Nodes 1,2,3 = 0	Node 1 = 600
	Node 2 = 140
	Node 3 = 64
<hr/>	<hr/>
Total = 0	Total = 804

Total Cost = 0 + 804 = 804 (No Replication)

*** Recommended File Placement for File 1: NODE 1, NODE 3 ***
TOTAL COST = 572

5) Perform Sensitivity Analysis (if desired)
(Algorithm Step 5)

6) ** Final Placement for File 1 **



7) Create Partitions/Files: Relation 2
(Algorithm Step 2)

Part A: Build Document Table (Sort by Node)

							made	ord		
	part	qty	made	qty	cord	part	to	to	qty	part
Node:Doc	no.	made	date	ord	date	desc	date	date	rem	price
1 CUST_ORD	20			20*	20*					
1 HIST_FILE	30			30	30					30
1 INVOICE	10			10	10					10
3 NEW_PART	2*					2*				2*
3 PART_SUP	10	10*	10*			10				10
3 CALC_INV	30	30	30	30	30		30*	30*	30*	
3 INV_REPT	1					1	1	1		

Part B: Analyze Frequency of Use, Determine "Owners"

qty_made
made_date: 30 queries, 10 updates - Node 3 (X)

qty_ord
cord_date: 40 queries, 20 updates - Node 1 (X)
30 queries, 0 updates - Node 3

part_desc: 11 queries, 2 updates - Node 3 (X)

made_to_date

ord_to_date

qty_rem: 1 query, 30 updates - Node 3 (X)

part_price: 40 queries, 0 updates - Node 1 (X)
10 queries, 2 updates - Node 3

Part C: Merge Attributes by Owner to form Partitions:

File 2 = (part_no, qty_ord, cord_date, part_price)

File 3 = (part_no, qty_made, made_date, part_desc,
made_to_date, ord_to_date, qty_rem)

Part D: Modify Partitions (if desired)

8) Determine File Placement: Relation 2, File Name= 2
(Algorithm Step 3)

Part A: Build Transaction Table: File 2 (algorithm step 4)
File 2 = (part_no, qty_ord, cord_date, part_price)

	Queries	Updates
1	40	20
2	0	0
3	40	2

Part B: Build Cost Tables

(Network Cost Table from Previous Step)

(Nodes)

	1	2	3
1	0	8	12
2	8	0	20
3	12	20	0

Queries

	1	2	3
1	0	320	480
2	0	0	0
3	480	800	0
	480	1120	480

Updates

	1	2	3
1	0	160	240
2	0	0	0
3	24	40	0
	24	200	240

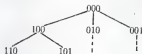
Total Costs: Node 1: $480 + 24 = 504$

Node 2: $1120 + 200 = 1320$

Node 3: $480 + 240 = 720$

*** Recommended File Placement at Node 1 **

- 9) Determine Profitability of Replication: File 2
(algorithm step 4)



- a) Node 1 and 2:

<u>Query Cost:</u>	<u>Update Cost:</u>
Node 1,2 = 0	Node 1 = 160
Node 3 = 480	Node 2 = 200
	Node 3 = 64
<hr/>	<hr/>
Total = 480	Total = 224

Total Cost = 480 + 224 = 704 (No replication)

- b) Node 1 and 3:

<u>Query Cost:</u>	<u>Update Cost:</u>
Node 1,3 = 0	Node 1 = 240
Node 2 = 0	Node 2 = 0
	Node 3 = 24
<hr/>	<hr/>
Total = 0	Total = 264

Total Cost = 0 + 264 = 264 ** Replicate **

- c) Continue to Branch, Nodes 1,2,3:

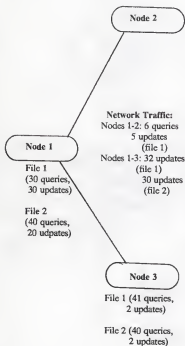
<u>Query Cost:</u>	<u>Update Cost:</u>
Nodes 1,2,3 = 0	Node 1 = 400
	Node 2 = 0
	Node 3 = 64
<hr/>	<hr/>
Total = 0	Total = 464

Total Cost = 0 + 464 = 464 (No Replication)

- * Recommended File Placement for File 2: NODE 1, NODE 3 *
TOTAL COST = 264

10) Perform Sensitivity Analysis (if desired)
(Algorithm Step 5)

11) ** Final Placement For File 1 and File 2 **



12) **Determine File Placement:** Relation 2, File Name= 3
(Algorithm Step 3)

Part A: Build Transaction Table: File 3 (algorithm step 4)

File 3 = (part_no, qty_made, made_date, part_desc,
made_to_date, ord_to_date, qty_rem)

	Queries	Updates
1	0	0
2	0	0
3	1	42

Part B: Build Cost Tables

(Network Cost Table from Previous Step)

(Nodes)

	1	2	3
1	0	8	12
2	8	0	20
3	12	20	0

Queries

	1	2	3
1	0	0	0
2	0	0	0
3	12	20	0
	12	20	0

Updates

	1	2	3
1	0	0	0
2	0	0	0
3	504	840	0
	504	840	0

Total Costs: Node 1: $12 + 504 = 516$

Node 2: $20 + 840 = 860$

Node 3: $0 + 0 = 0$

***** Recommended File Placement at Node 3 ****

13) Determine Profitability of Replication: File 3
(algorithm step 4)

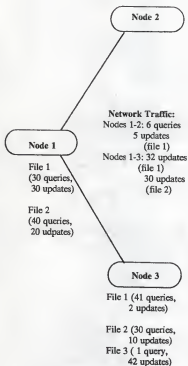
**** Cost is zero; replication not profitable ****

*** Recommended File Placement for File 3; NODE 3 ***
TOTAL COST = 0

14) Perform sensitivity analysis (If desired)
(Algorithm Step 5)

15) ** FINAL OUTPUT OF DIST_DATA TOOL **

File 1: at Nodes 1 and 3	Cost = 572
File 2: at Nodes 1 and 3	Cost = 264
File 3: at Nodes 2	Cost = 0
Total Cost	836



APPENDIX IV CODE LISTING: DIST_DATA.C

```

1  /.....
2  .....
3  .....
4  DIST_DATA_PROGRAM:
5  .....
6  The dist_data program was designed to assist the database designer in determining
7  correct file partitioning and placement in a distributed database environment.
8  The program takes as input the following files:
9  .....
10 - Bernfile: This is a file produced by execution of the 'bern2' program, which
11             produces 3NF Relations
12 - Docfile: This is a file produced by the 'dochandler' program that takes customer
13            documents and converts them to functional dependencies
14 .....
15 The program produces as output the following information:
16 .....
17 - Recommended Partitioning of the 3NF Relations
18 - Recommended File Placement/Replication of these partitions onto nodes in a distributed
19   network
20 .....
21 The main() function controls dist_data processing. It solicits information from the
22 user regarding network nodes, organizational structure and document usage per organization
23 once this data is gathered it begins processing to determine file partitioning
24 and file placement. The program consists of the following functions
25 .....
26 build_bern(): Builds the bern2 (3NF) relation table
27 build_rel(): Builds a relation table which indicates attribute usage on a per-node basis
28 find_test_cost(): Determines minimum cost connection between nodes on the network
29 build_parts(): Recommends file partitioning of the 3NF relations
30 pr_parts(): Prints resulting partitions created
31 file_placement(): Determines recommended file placement for each partition created
32                  (based on minimum cost for update/query usage)
33 branch_bound(): Performs branch and bound search to determine the profitability of
34                  replication
35 find_cost(): Finds least cost path between 2 nodes in a network
36 find_new(): Finds new cost associated with file replication
37 pr_reltbl(): Prints out contents of relation table
38 rebuild_cost(): Rebuilds the cost tables if sensitivity analysis has been requested
39                  {i.e., determine if replication is desirable if a 92.102, etc.
40                  communication/increase/decrease is incurred}
41 .....
42 .....
43 .....
44 #include <stdio.h>
45 #include <string.h>
46 #include <conio.h>
47 #include <fcntl.h>
48 #include <sys/types.h>

```

```

49 Int NO_NODES=0, NO_ATTR=0, NO_OOCS=0, NO_ORGS=0;
50 Int NO_RELS=0, NO_ENTS=0, holdgrp=0;
51
52 main(argc, argv)
53 Int argc;
54 char *argv[];
55 {
56     FILE *docptr, *bernptr, *fopen();
57     char *calloc();
58
59     Int c=0, l=0, j=0, k=0, chkcmp=0, tpcost=0, tpepnode=0, relnum=0;
60     Int search_node=0, least_cost=0, con_node = 0;
61     Int relentry=0, no_parts=0, partnumb=0;
62     Int repl_skip=0, tpe2_skip=0, operation = 0;
63     Int flag=FALSE, flag2=FALSE, percent=0;
64     char sign;
65     char ans[4];
66     char tpepnode[ORG_SZ];
67     char tpepnode2[ORG_SZ];
68     char tpepnode3[ORG_SZ];
69     if(argc !=3) /* Check for Valid No. of Args */
70     {
71         fprintf(stderr, "Usage: dist_data <docfile><bernpfile>\n");
72         exit(1);
73     }
74     else if((docptr=fopen(++argv, "r")) == NULL)
75     {
76         fprintf(stderr, "dist_data: cannot open %s\n", *argv);
77         exit(1);
78     }
79     else if((bernptr=fopen(++argv, "r")) == NULL)
80     {
81         fprintf(stderr, "dist_data: cannot open %s\n", *argv);
82         exit(1);
83     }
84     /* GET TOTAL NO. OF OOCS/ATTRIBUTES FROM OOCFILE */
85     fscanf(docptr, "%d%d", &NO_OOCS, &NO_ATTR);
86
87     /* DETERMINE NUMBER OF BERN2 RELATIONS */
88     while ((c=fgetc(bernptr)) != EOF)
89     {
90         if(c=='(')
91             NO_RELS++;
92     }
93
94
95
96

```



```

97 }
98 fseek(berptr,0L,0);
99
100 /* INITIALIZE BERTABLE */
101
102 for(i=0; i<NO_RELS; i++)
103 {
104     for(j=0; j<NO_ATTR; j++)
105     {
106         strcpy(berstable[i].attributes[j].attr, "");
107         berstable[i].attributes[j].key = FALSE;
108     }
109 }
110
111 /* BUILD BERN TABLE (holds bern2 relations) */
112 build_bern(berptr);
113
114 /* QUERY USER */
115
116 printf("Enter No. Nodes on Network: ");
117 scanf("%d",&NO_NODES);
118 printf("\nEnter No. of Organizations: ");
119 scanf("%d",&NO_ORGS);
120
121 /* INITIALIZE COST TABLE */
122 for(i=0; i<NO_NODES; i++)
123 {
124     for(j=0; j<NO_NODES; j++)
125         costable[i][j] = 0;
126 }
127
128 /* INITIALIZE PORTION OF RELATION TABLE */
129 for(i=0; i< MAX_ENTRY; i++)
130 {
131     reltable[i].nodenum = 0;
132     strcpy(reltable[i].orgname, "");
133     strcpy(reltable[i].docname, "");
134 }
135
136 /* QUERY USER */
137
138 for(i=0; i<NO_ORGS; i++)
139 {
140     strcpy(tmpdoc, "");
141     strcpy(tmporg, "");
142
143
144

```

```

145 printf("\nEnter Org. Name: ");
146 scanf("%s",&tporg);
147
148
149 printf("\n%s on node?: ",&tporg);
150 scanf("%d",&tnode);
151
152 printf("\nEnter Interconnect Node(s) & Cost(s): Enter each node/cost on separate line.\n");
153 for(j=0; j< NO_NODES; j++)
154 {
155     con_node=0;
156     scanf("%d",&con_node);
157
158     if(con_node != 0)
159     {
160         scanf("%d",&tmpcost);
161         costable[tnode-1][con_node-1] = tmpcost;
162         costable[con_node-1][tnode-1] = tmpcost;
163         tmpcost=0;
164     }
165     else
166     {
167         break;
168     }
169 } /* End of Nodes for Organization */
170
171
172 printf("\nEnter Documents Owned: Enter each document on separate line.\n");
173 printf("enter 'quit' to complete\n");
174 for(j=0; j<NO_DOCS; j++)
175 {
176     strcpy(tmpdoc,"");
177     scanf("%s",tmpdoc);
178     if((strcmp(tmpdoc,"quit")==0) || (strcmp(tmpdoc,"QUIT")==0))
179     {
180         break;
181     }
182     else
183     {
184         reliable[reentry].nodenum = tnode;
185         strcpy(reliable[reentry].orgname,tporg);
186         strcpy(reliable[reentry].docname,tmpdoc);
187         reentry++;
188     }
189 } /* End of Organization Documents */
190
191 } /* End of Processing Organizations */
192 NO_ENTS=reentry;

```

```

193 /* CHECK COST TABLE */
194
195 for(i=0; i< NO_NODES; i++)
196 {
197     for(j=0; j<NO_NODES; j++)
198     {
199         if(costable[i][j] == 0)
200             if (i != j)
201             {
202                 search_node = j;
203                 tmpcost=find_least_cost(i,j,search_node,least_cost);
204                 costable[i][j] = tmpcost;
205                 costable[j][i] = tmpcost;
206                 least_cost = 0;
207             }
208         }
209     }
210     else
211     {
212         costable[i][j] = 0;
213         costable[j][i] = 0;
214     }
215 }
216 } /* End of for Loop */
217
218 } /* End of for Loop */
219
220
221 /* DIST_DATA PROCESSING */
222
223 for(i=0; i< NO_RELS; i++)
224 {
225     /* INITIALIZE RELATION TABLE */
226     for (j=0; j<MAX_ENTRY; j++)
227     {
228         retable[j].use_rel=FALSE;
229         retable[j].total_query=0;
230         retable[j].total_update=0;
231         for(k=0;k<MAX_ATTRS; k++)
232         {
233             retable[j].freq[k].key=FALSE;
234             retable[j].freq[k].owner=FALSE;
235             retable[j].freq[k].part_num=0;
236             retable[j].freq[k].frequency=0;
237         }
238     }
239 }
240

```

```

241 build_rel(docptr,i);
242 no_parts=build_parts(1);
243
244 for(k=0; k<MAX_PARTS; k++)
245 {
246     if(parts_created[k] != -1)
247     {
248         partnum=parts_created[k];
249         file_placement(1,partnum);
250         printf("\nPERFORMING REPLICATION ANALYSIS (Branch & Bound)...\n");
251         sleep(2);
252         branch_bound();
253     }
254     /* SEE IF SENSITIVITY ANALYSIS DESIRED */
255     while (flag == FALSE)
256     {
257         printf("\nSensitivity Analysis Desired ? ");
258         scanf("%s",&ans);
259         if((ans[0] == 'y') || (ans[0] == 'Y'))
260         {
261             flag=TRUE;
262         }
263         /* RE-INITIALIZE SOLUTION ARRAY */
264         for (j=0; j< MAX_NODES; j++)
265         {
266             if(j != evnode - 1)
267             {
268                 solution_array[j] = 0;
269             }
270             printf("\ninput will be taken as a percentage increase/decrease\n");
271             while (flag2 == FALSE)
272             {
273                 printf("\nEnter '+' for Increase/ '-' for Decrease: ");
274                 scanf("%s",&sign);
275                 if((sign == '+') || (sign == '-'))
276                 {
277                     flag2=TRUE;
278                 }
279                 else
280                 {
281                     printf("\nMust indicate '+' or '-' (increase or decrease)\n");
282                 }
283                 printf("\nEnter Cost Factor (1.0., 5, 50, etc): ");
284                 scanf("%d",&percent);
285                 rebuild_cost(sign,percent);
286                 branch_bound();
287             }
288         }

```

```

289     else if((ans[o] == 'n') || (ans[o] == 'N'))
290         flag=TRUE;
291     else
292         printf("Answer must be Yes or No\n");
293     }
294     flag2=FALSE;
295     }
296     flag=FALSE;
297     }
298
299     /* RESET FILE POINTER FIRST DOCUMENT (strip off total
300     document and attribute counters) */
301
302     fseek(docptr,0L,0);
303     fscanf(docptr,"%d%d",&tmp1_skip,&tmp2_skip);
304
305     } /* End of Relations */
306
307     /* COMMENT OUT PRINT OF THESE FILES
308     pr_costbl();
309     pr_berntbl();
310     */
311     } /* End of Main */
312
313
314     /*****
315     * BUILD_BERN():
316     *
317     * Processing: Parses the bern2 file to determine 3MF relations, and to flag
318     * which attributes are key fields in the relation
319     *
320     * Input: File pointer to output file produced by bern2 (3MF Relations)
321     *
322     * Functions Called: (None)
323     *
324     * Output/ Return: Populates a structure which contains the bern2 relations for
325     * later use in partitioning functions.
326     *
327     *****/
328
329     build_bern(fp)
330     char *fp;
331     {
332         int relnum=0, attrct=0, i=0;
333         int keyflag = TRUE;
334
335

```

```

337 int validchar = FALSE;
338
339 int holdchar;
340 holdchar = getch(fp); /* Get rid of leading '\n' */
341 while (holdchar != EOF)
342 {
343     holdchar = getch(fp);
344     switch(holdchar)
345     {
346         case '\n':
347             case '\t':
348                 case '>':
349                     validchar=FALSE;
350                     break;
351             case '{':
352                 keyflag = TRUE;
353                 attrnum[relnum] = attrct;
354                 relnum++;
355                 attrct=0; i=0;
356                 validchar=FALSE;
357                 break;
358             case ',':
359                 keyflag=FALSE;
360                 validchar=FALSE;
361                 break;
362             case ':':
363                 if(!validchar==TRUE) && (keyflag==TRUE))
364                 {
365                     berntable[relnum].attributes[attrct].key = TRUE;
366                     berntable[relnum].attributes[attrct].attrs[i] = '\0';
367                     attrct++;
368                     i=0;
369                 }
370                 if(validchar==TRUE)
371                 {
372                     berntable[relnum].attributes[attrct].attrs[i] = '\0';
373                     attrct++;
374                     i=0;
375                 }
376                 validchar=FALSE;
377                 break;
378             case EOF:
379                 if(!validchar==TRUE) && (keyflag==TRUE))
380                 {
381                     berntable[relnum].attributes[attrct].key = TRUE;
382                     berntable[relnum].attributes[attrct].attrs[i] = '\0';
383                     attrct++;
384                     i=0;
385                 }
386                 if(validchar==TRUE)
387                 {
388                     berntable[relnum].attributes[attrct].attrs[i] = '\0';
389                     attrct++;
390                     i=0;
391                 }
392             }
393     }

```

```

385     }
386     attrnum[relnum] = attrct;
387     break;
388     default:
389         berntable[relnum].attributes[attrct].attr[i] = holdchar;
390         validchar=TRUE;
391         i++;
392     } /* End of "switch" */
393 } /* End of "while" */
394 }
395
396 /*****
397  * FIND_LEAST_COST():
398  *
399  * Processing: Recursive function that traces path connections between nodes
400  * on the network, to determine the "least cost" path between
401  * any given pair of nodes (the "least data processing requires
402  * a fully connected network (virtual) in order to have a cost
403  * associated between any two nodes". This is needed for proper
404  * file placement based on minimum cost).
405  *
406  * Input: Index values (i,j) of network cost table, search node which
407  * which is the node representing an "open" connection and
408  * min_cost which is the placeholder for the minimum cost found
409  * so far in the search
410  *
411  * Functions: find_least_cost() (Recursive function)
412  *
413  * Called:
414  *
415  * Output:
416  *
417  * Return: Returns the least cost connection between a given pair of nodes
418  *
419  * int i,j,search_node,mincost;
420  *
421  * {
422  *     int k, new_cost,add_cost = 0;
423  *     int NO_NODES = 5;
424  *     for(k=0; k<NO_NODES;k++)
425  *     {
426  *         if(costable[search_node][k] !=0 && (k != j))
427  *         {
428  *             add_cost = add_cost + costable[search_node][k];
429  *             if(costable[i][k] !=0)
430  *             {
431  *                 add_cost=add_cost + costable[i][k];
432  *

```

```

433     if(mincost == 0)
434         mincost = add_cost;
435     else if(add_cost < mincost)
436         mincost = add_cost;
437     else
438         add_cost=0;
439 }
440 {
441     search_node = k;
442     find_least_cost(i,j,search_node,add_cost,mincost);
443 }
444 }
445 }
446 return(mincost);
447 }
448
449 /******
450 * BUILD_REL():
451 * Processing: Parses the bern2 file to determine 3NF relations, and to flag
452 *              which attributes are key fields in the relation.
453 *
454 * Input:       File pointer to output file produced by bern2 (3NF Relations)
455 *
456 * Functions   Called: (None)
457 *
458 * Output:     Populates a structure which contains the bern2 relations for
459 *              later use in partitioning functions
460 *
461 * Return:     later use in partitioning functions
462 *
463 *
464 *
465 *
466 *
467 *
468 *
469 *
470 *
471 *
472 *
473 *
474 *
475 *
476 *
477 *
478 *
479 *
480 */
481
482 build_rel(fp,relno)
483 char *fp;
484 int relno;
485 {
486     int validchar=FALSE,attowner=FALSE,keyflag=FALSE;
487     char tmpoccl[DOC_SZ];
488     char tmpatt[DOC_SZ];
489     char tmpkey[DOC_SZ];
490     int check=0,duplicate=0;
491     int tmpfreq=0,attret=0,tblent=0,attrsz=0,i=0;
492     strcpy(tmpatt,"");
493     dumpchar = fgetc(fp);

```



```

481 while (check != EOF)
482 {
483     check=fgetc(fp);
484     if(check == '\n')
485     {
486         fscanf(fp,"%s",tmpdoc);
487     }
488     /* SKIP THE "MISC" IN DOCUMENT LIKE LOCATION.INPUT, ETC */
489     for (i=0; i< 4; i++)
490     {
491         fscanf(fp,"%s",tmpjunk);
492     }
493     fscanf(fp,"%d",&tmpfreq);
494     check=fgetc(fp);
495     if(check != '\n')
496     {
497         else
498         {
499             switch(check)
500             {
501                 case '\t':
502                 case '\n':
503                 case ' ':
504                 case '+':
505                 {
506                     tmpattr[attrsz] = '\0';
507                     for(attrct=0; attrct< (attrnum[relno]+1); attrct++)
508                     {
509                         if(strncmp(tmpattr,berntable[relno].attributes[attrct].attrs) == 0)
510                         {
511                             if(berntable[relno].attributes[attrct].key == TRUE)
512                                 keyflag=TRUE;
513                             for(tblent=0; tblent < NO_ENTS+1; tblent++)
514                             {
515                                 if(strncmp(tmpdoc,reltable[tblent].docname) == 0)
516                                 {
517                                     reltable[tblent].use_rel=TRUE;
518                                     reltable[tblent].freq[attrct].frequency=tmpfreq;
519                                     if(attrowner == TRUE)
520                                     {
521                                         reltable[tblent].freq[attrct].owner = TRUE;
522                                         if(keyflag==TRUE)
523                                             reltable[tblent].freq[attrct].key=TRUE;
524                                     }
525                                 }
526                             }
527                         }
528                     }
529                     /* Finish search of table for one attr */
530                     /* Finish "if" statement */

```

```

531 } /* Finish "for" statement */
532 validchar=FALSE;
533 if (validchar)
534   strcpy(tmpattr,"");
535   attrname = FALSE;
536   keyflag=FALSE;
537   }
538   break;
539   case 'X':
540     attrowner=TRUE;
541     break;
542   default:
543     tmpattr[attrsz] = check;
544     attrsz++;
545     validchar=TRUE;
546     break;
547   }
548   } /* End of switch statement */
549   } /* End of "else" statement */
550   } /* End of "while" statement */
551 }
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578

```

BUILD_PARTS():

Processing: Analyzes each attribute in the relation table to determine how often this attribute is used at each node. It then assigns the attribute to the node where it is most frequently used. After all attributes are processed, the attributes are grouped together by frequency of use (i.e., all attributes which are used most frequently at node 'X' will be grouped together, thus forming partitions. After the partitioning has been completed, this is function assigns the used to each additional partition, allowing if needed (the number of partitions is determined by the number of new partitions, above an attribute from one partition to another or MERGE two partitions)

Input: Relation number which is to be processed, or analyzed for partitioning

Functions

Called: pr_parts()

Output/

Return: Populates the relation table with a partition number for each attribute and returns the number of partitions created.

```

579 build_parts(rerno)
580 int rerno;
581 {
582     struct ATTCT
583     {
584         int t_gry;
585         int t_upd;
586     }attct[MAX_NODES];
587
588     int count=0, tblend=0, NODE_NUM=0, keypno=0, chk_freq=0, l=0, most_freq=0;
589     int node=0, m_upd=0, num_parts=0;
590     int flag=FALSE, modflag=FALSE, found_newnumb=FALSE, NEW_PART=TRUE;
591     int chkflag=FALSE, chkflag2=FALSE, ansflag=0, att_number=0, move_part=0;
592     int part_no=0, part1=0, part2=0, operation = 0;
593     char ans[4], attribute[ATR_52];
594
595     keypno=NO_NODES+1;
596
597     /* INITIALIZE PARTITION ARRAY */
598     for (i=0; i< MAX_PARTS; i++)
599         parts_created[i] = -1;
600
601     /* DISPLAY RESULTS TO USER AS ANALYSIS IS BEING PERFORMED */
602     printf("\n\nANALYZING ATTRIBUTE USAGE FOR RELATION %d ....\n", (rerno+1));
603
604     /* LOOP THROUGH ALL ATTRIBUTES TO ANALYZE */
605     for(count=0; count< (attrnum[rerno+1]); count++)
606     {
607         /* INITIALIZE ATTCT */
608         for(i=0; i<MAX_NODES;i++)
609         {
610             attct[i].t_gry=0;
611             attct[i].t_upd=0;
612         }
613
614         /* IF ATTRIBUTE IS A KEY -- IF IT IS NOT INCLUDED IN ANALYSIS */
615         /* (as a key will appear in all partitions created) */
616         if(reftable[0].freq[count].key != TRUE)
617         {

```

```

627 /* LOOP THROUGH TABLE ENTRIES FOR FREQ. OF USE */
628
629 for(tblent=0; tblent<ND_ENTS+1; tblent++)
630 {
631     NODE_NUM=reliable[tblent].nodenum-1;
632     if(reliable[tblent].freq[count].owner == TRUE)
633         attct[NODE_NUM].t_upd = attct[NODE_NUM].t_upd + reliable[tblent].freq[count].frequency;
634     else
635         attct[NODE_NUM].t_qry = attct[NODE_NUM].t_qry + reliable[tblent].freq[count].frequency;
636 }
637
638 /* DETERMINE "OWNER" (Node most frequently accessed) */
639 most_freq=0; m_upd=0; node=0;
640
641 /* PRINT FOR USER */
642
643 if(reliable[0].freq[count].key != TRUE)
644     printf("\ndebernabi[reld].attributes[count].attrs);
645     for(i=0; i<ND_NODES;i++)
646     {
647         if(reliable[0].freq[count].key != TRUE)
648             printf("\n%d queries/12d updates (Node %d)", attct[i].t_qry, attct[i].t_upd, i+1);
649         chk_freq = attct[i].t_upd + attct[i].t_qry;
650         if(chk_freq > most_freq)
651         {
652             most_freq = chk_freq;
653             m_upd = attct[i].t_upd;
654             node = i;
655         }
656         else if (chk_freq == most_freq)
657         {
658             if (attct[i].t_upd > m_upd)
659             {
660                 most_freq = chk_freq;
661                 m_upd = attct[i].t_upd;
662                 node = i;
663             }
664         }
665     }
666
667 /* POPULATE ATTRIBUTE WITH PARTITION NUMBER
668 /* (Partition number equals owner node no.
669
670 for(tblent=0; tblent<ND_ENTS+1; tblent++)
671 {
672     if(reliable[0].freq[count].key != TRUE)
673 
```

```

675 {
676     reltable[tblent].freq[count].part_num=node+1;
677 }
678 else
679 {
680     reltable[tblent].freq[count].part_num = keypno;
681 }
682 }
683 /* KEEP TRACK OF NUMBER OF PARTITIONS CREATED */
684 NEW_PART=TRUE;
685 for(i=0; i< NO_NODES; i++)
686 {
687     if(node+1 == parts_created[i])
688         NEW_PART=FALSE;
689 }
690 if(NEW_PART == TRUE)
691 {
692     parts_created[numb_parts] = node+1;
693     numb_parts++;
694 }
695 }
696 /* PRINTOUT RESULTING PARTITIONS CREATED */
697 pr_parts(numb_parts,relno);
698 /* SEE IF DISPLAY OF PARTITIONING IS DESIRED */
699 while (flag == FALSE)
700 {
701     printf("\ndisplay of Partitioning Analysis Desired ? ");
702     scanf("%s",&ans);
703     if((ans[0] == 'y') || (ans[0] == 'Y'))
704     {
705         flag=TRUE;
706         pr_reltbl(relno);
707     }
708     else if((ans[0] == 'n') || (ans[0] == 'N'))
709         flag=TRUE;
710     else
711         printf("Answer must be Yes or No\n");
712 }
713 /* SEE IF MODIFICATIONS TO PARTITIONING ARE DESIRED */
714 flag = FALSE;
715
716
717
718
719
720
721
722

```

```

723 while (flag == FALSE)
724 {
725     printf("\nModification of Partitioning Desired ? ");
726     scanf("%s", ans);
727     if((ans[0] == 'y') || (ans[0] == 'Y'))
728     {
729         flag=TRUE;
730         modflag=TRUE;
731     }
732     else if((ans[0] == 'n') || (ans[0] == 'N'))
733     flag=TRUE;
734     else
735         printf("\nAnswer must be Yes or No\n");
736 }
737
738 while(modflag==TRUE)
739 {
740     operation = 0;
741     printf("Operations Available:\n");
742     printf("1 = MERGE two partitions\n");
743     printf("2 = MOVE an attribute to a different partition\n");
744     printf("3 = CREATE a new partition\n");
745     printf("4 = quit (done making modifications)\n");
746     printf("\nEnter Number of Desired Operation: ");
747     scanf("%d", &operation);
748     switch(operation)
749     {
750     case 1:
751         ansflag = FALSE, chkflag1 = FALSE, chkflag2 = FALSE;
752         parts_O, part2 = 0;
753         if(numb_parts == 1)
754         {
755             printf("\nInvalid operation -- only 1 partition exists\n\n");
756             operation = 0;
757         }
758         else
759         {
760             while (ansflag == FALSE)
761             {
762                 printf("This operation will merge two existing partitions into one\n");
763                 printf("-in the form of MERGE PARTITION ____ WITH ____\n\n");
764                 printf("Enter File Partition Numbers: ");
765                 scanf("%d%d", &part1, &part2);
766                 for(i=0; i<numb_parts; i++)
767                 {
768                     if(part1 == parts_created[i])
769                         chkflag1=TRUE;
770                     if(part2 == parts_created[i])

```

```

771         }
772         chkflag2=TRUE;
773     }
774     if((chkflag1 == TRUE) && (chkflag2 == TRUE))
775         ansflag = TRUE;
776     else
777     {
778         ansflag = FALSE;
779         if(chkflag1 == FALSE)
780             printf("First Partition Entered is Invalid partition number\n");
781         else
782             printf("Second Partition Entered is Invalid partition number\n");
783     }
784 }
785 }
786 break;
787
788 case 2:
789
790     ansflag = FALSE, chkflag1 = FALSE, chkflag2 = FALSE, move_part=0;
791     while (ansflag == FALSE)
792     {
793         printf("This Operation will move an attribute from one partition to another\n");
794         printf("In the form of MOVE ATTRIBUTE ____ TO PARTITION ____\n");
795         printf("Enter Attribute Name: ");
796         scanf("%s", attribute);
797         printf("Enter Partition Number: ");
798         scanf("%d", &move_part);
799     }
800     /* VERIFY A VALID ATTRIBUTE WAS ENTERED */
801     for(i=0; i< attrnum[reino] +1; i++)
802     {
803         if(strcmp(attribute, barntable[reino].attributes[i].attrs) == 0)
804         {
805             chkflag1 = TRUE;
806             att_number = i;
807         }
808     }
809     /* VERIFY VALID PARTITION NUMBER WAS ENTERED */
810     for(i=0; i<numb_parts; i++)
811     {
812         if(move_part == parts_created[i])
813         {
814             chkflag2=TRUE;
815         }
816     }

```

```

819     }
820     if((chkflag1 == TRUE) && (chkflag2 == TRUE))
821     {
822         ansflag = TRUE;
823     }
824     else
825     {
826         ansflag = FALSE;
827         if(chkflag1 == FALSE)
828             printf("Invalid attribute name entered\n");
829         else
830             if(chkflag2 == FALSE)
831                 printf("Partition Entered is Invalid partition number\n");
832     }
833     break;
834 }
835 case 3:
836     printf("This Operation creates a new file partition\n");
837     printf("This partition can then be populated by a series of MOVE operations\n\n");
838     while(found_newnumb==FALSE)
839     {
840         npart_no++;
841         found_num=TRUE;
842         for(i=0; i<numb_parts;i++)
843         {
844             if((npart_no == NO_PARTS +1) || (npart_no == parts_created[i]))
845                 found_newnumb=FALSE;
846         }
847     }
848     parts_created[numb_parts] = npart_no;
849     printf("*** PARTITION NUMBER CREATED = %d ***\n",npart_no);
850     sleep(2);
851     ansflag=TRUE;
852     numb_parts++;
853     break;
854 }
855 case 4:
856     modflag = FALSE;
857     ansflag = FALSE;
858     break;
859 default:
860     printf("Operation must be Numbers 1-4 \n");
861 }
862 } /* End of "switch" statement */
863
864 /* CHECK OPERATION REQUESTED AND MAKE APPROPRIATE CHANGES TO RELIABLE */
865
866

```



```

867 if((operation == 1) || (operation == 2))
868 {
869     for (tblent=0; tblent < NO_ENTS +1; tblent++)
870     {
871         if(operation == 2)
872             reliable[tblent].freq[att_number].part_num = move_part;
873         if(operation == 1)
874         {
875             for(count=0; count < (attrnum[relno] +1); count++)
876             {
877                 if(reliable[tblent].freq[count].part_num == part1)
878                 {
879                     reliable[tblent].freq[count].part_num = part2;
880                 }
881             }
882         }
883     }
884 }
885 if(operation == 1)
886 {
887     for(i=0; i< numb_parts; i++)
888     {
889         if(parts_created[i] == part1)
890         {
891             parts_created[i] = -1;
892             numb_parts--;
893         }
894     }
895 }
896 if(ansflag == TRUE)
897 {
898     printf("\nFILE PARTITIONS AFTER MODIFICATIONS\n");
899     pr_parts(numb_parts, relno);
900 }
901 /* End of "while" statement */
902 return(numb_parts);
903 }
904 }
905 }
906 }
907 }
908 }
909 }
910 }
911 }
912 }
913 }
914 }

```

```

915 * Functions (None)
916 *
917 * Called:
918 *
919 * Output:
920 * Return: simply prints to output
921 *
922 *
923 *
924 *
925 *
926 *
927 *
928 *
929 *
930 *
931 *
932 *
933 *
934 *
935 *
936 *
937 *
938 *
939 *
940 *
941 *
942 *
943 *
944 *
945 *
946 *
947 *
948 *
949 *
950 *
951 *
952 *
953 *
954 *
955 *
956 *
957 *
958 *
959 *
960 *
961 *
962 *

```

```

pr_parts(no_parts,relnumb)
{
    int no_parts,relnumb;

    int attrno, i;
    unsigned int part_pr;
    int PART_FOUND = FALSE;

    printf("\n*****\n");
    printf("FILE PARTITION(S) CREATED FOR RELNO %d\n",relnumb);
    printf("(File Partition Names Assigned According To 'Node/Owner')\n");

    for (i=0; i< MAX_PARTS; i++)
    {
        part_pr = parts_created[i];
        if(part_pr != -1)
        {
            printf("\nFILE PART %d: ",part_pr);
            for(attrno=0; attrno < (attrnum[relnumb] +1); attrno++)
            {
                if(releable[0].freq[attrno].part_num==NO_NODES + 1)
                {
                    printf("%s > ",berntable[relnumb].attributes[attrno].attrs);
                }
                else if(part_pr == releable[0].freq[attrno].part_num)
                {
                    printf("%s",berntable[relnumb].attributes[attrno].attrs);
                }
                printf("\n\n*****\n\n\n");
            }
        }
    }

    FILE_PLACEMENT();

    Processing: Works on one file partition at a time, to build transaction table
    (aggregate update/query access to this partition from each node),
    builds the associated transaction cost tables (query and update

```

```

963 * costs to this file, if placed on any given node), and determines
964 * recommended file placement based on least-cost.
965 *
966 * Input: The Relation number and partition number to be analyzed for file placement
967 *
968 * Functions (None)
969 * Called:
970 *
971 * Output/
972 * Return: Sets file placement bit for this partition in global solution_array
973 * (returns nothing)
974 * ...../
975 *
976 * file_placement[relnum,part_no]
977 * int relnum,part_no;
978 * {
979 *
980 *
981 * int i=0,j=0,tblent=0,count=0,svfreq=0,upd_trans=FALSE; int calc_total=0,calc_upd=0,calc_qlry=0.
982 * svnode=0, svmin_cost=0;
983 *
984 * /* INITIALIZE TRANSACTION/COST TABLES/ SOLUTION ARRAY */
985 *
986 * for (i=0; i<MAX_NODES; i++)
987 * {
988 *   trans_tbl[i].t_qlry=0;
989 *   trans_tbl[i].t_upd=0;
990 *   solution_array[i] =0;
991 *
992 *   for(j=0; j<MAX_NODES; j++)
993 *   {
994 *     qry_cost[i][j] = 0;
995 *     upd_cost[i][j] = 0;
996 *
997 *   }
998 * }
999 *
1000 * /* SEARCH RELTABLE TO FIND TOTAL UPD AND QUERY FOR EACH NODE */
1001 * /* {Keys are included in total access }
1002 *
1003 * for(tblent=0; tblent < NO_ENTS +1; tblent++)
1004 * {
1005 *   svnode = rellable[tblent].noderun - 1;
1006 *   for(count=0; count < attrnum[relnum] +1; count++)
1007 *   {
1008 *     if(rellable[tblent].freq[count].part_num==part_no)
1009 *     {
1010 *       if(rellable[tblent].freq[count].owner == TRUE)

```

```

1011         if(svfreq < reitablen[tblent].freq[count].frequency)
1012             svfreq = reitablen[tblent].freq[count].frequency;
1013     }
1014 }
1015
1016 if(UPD_TRANS == TRUE)
1017 {
1018     trans_tbl[svnode].t_upd = trans_tbl[svnode].t_upd + svfreq;
1019     UPD_TRANS=FALSE;
1020 }
1021 else
1022     trans_tbl[svnode].t_qry = trans_tbl[svnode].t_qry + svfreq;
1023     svfreq= 0;
1024 }
1025
1026 /* PRINT TRANSACTION TABLE */
1027 printf("\n\nTRANSACTION TABLE FOR RELND:%d FILE PARTITION: %d\n\n",relnum+1,part_no);
1028 printf("\n\nQUERY\UPDATE\n");
1029 for(i=0; i<ND_NODES; i++)
1030     printf("\t %d\t %d\t %d\n",i+1,trans_tbl[i].t_qry,trans_tbl[i].t_upd);
1031
1032 /* BUILD COST TABLES FOR QUERIES AND UPDATES */
1033 for(i=0; i<ND_NODES; i++)
1034 {
1035     for(j=0; j<ND_NODES; j++)
1036     {
1037         qry_cost[i][j] = trans_tbl[i].t_qry * costtab[i][j];
1038         upd_cost[i][j] = trans_tbl[i].t_upd * costtab[i][j];
1039     }
1040 }
1041
1042 /* DETERMINE RECOMMENDED FILE PLACEMENT FOR FILE PARTITION */
1043 printf("\n\nTOTAL COSTS ASSOCIATED WITH FILE PLACEMENT:\n\n");
1044 for(j=0; j<ND_NODES; j++)
1045 {
1046     for(i=0; i<ND_NODES; i++)
1047     {
1048         calc_qry = qry_cost[i][j] + calc_qry;
1049         calc_upd = upd_cost[i][j] + calc_upd;
1050     }
1051     calc_total = calc_qry + calc_upd;
1052     printf("NODE %d: %d Query + %d Update = %d Total Cost\n",j+1,calc_qry,calc_upd,calc_total);

```

```

1059 /* CHECK TO SEE IF MINIMUM COST */
1060 if(j==0)
1061 {
1062     svnode = j+1;
1063     min_cost = calc_total;
1064 }
1065 else if(calc_total < min_cost)
1066 {
1067     min_cost = calc_total;
1068     svnode = j+1;
1069 }
1070 calc_qry=0, calc_upd=0, calc_total=0;
1071 }
1072 svmin_cost = min_cost;
1073 printf("\n***** RECOMMENDED FILE PLACEMENT AT NODE %d. TOTAL COST= %d ***\n", svnode, min_cost);
1074 sleep(5);
1075 /* SET FILE PLACEMENT BIT FOR THIS NODE */
1076 solution_array[svnode-1] = 1;
1077 }
1078 }
1079 }
1080 }
1081 }
1082 }
1083 }
1084 }
1085 }
1086 }
1087 }
1088 }
1089 }
1090 }
1091 }
1092 }
1093 }
1094 }
1095 }
1096 }
1097 }
1098 }
1099 }
1100 }
1101 }
1102 }
1103 }
1104 }
1105 }
1106 }

```

Processing: Analyzes profitability of replication, by performing a branch and bound search. After a single file placement has been determined, this function will see if 2 copies of the file are profitable, 3 copies, etc. If at any point, the cost of replication is larger than the minimum cost, the "branching" stops for this node. Should a total cost of placing copies of the file at multiple nodes be less than the previous minimum; replication is deemed profitable. The results of the replication analysis (and associated costs) are displayed to the user.

Input: None

Functions Called: find_cost()

Output/Return: Updates global solution array to set new file placement bit(s) if replication was found to be profitable.

```

1107 branch_bound()
1108 {
1109     int keep_cost[MAX_NODES]; /* Keep intermediate cost figures */
1110     int i=0, node=0, level=1, new_cost=0, mincost_node=0, CONTINUE_BRANCH=TRUE;
1111     /* INITIALIZE KEEP_COST ARRAY */
1112     for(node=0; node < MAX_NODES; node++)
1113         keep_cost[node] = 0;
1114     /* MAIN ENTRY FOR BRANCH AND BOUND SEARCH */
1115     while((CONTINUE_BRANCH == TRUE) && (level < NO_NODES))
1116     {
1117         /* RESET BRANCH FLAG FOR NEXT PASS */
1118         /* MOVE DOWN 1 LEVEL
1119         CONTINUE_BRANCH= FALSE;
1120         level++;
1121         for(node=0; node < NO_NODES; node++)
1122         {
1123             if(solution_array[node] != 1)
1124             {
1125                 new_cost=find_cost(node);
1126                 if(((level==2)&&(new_cost<svmin_cost))|((level==2)&&(new_cost<min_cost)))
1127                 {
1128                     min_cost = new_cost;
1129                     mincost_node= node;
1130                     CONTINUE_BRANCH=TRUE;
1131                 }
1132             }
1133         }
1134         if(CONTINUE_BRANCH == TRUE)
1135             solution_array[mincost_node] = 1;
1136     }
1137     printf("\n*** AFTER REPLICATION ANALYSIS, RECOMMENDED FILE PLACEMENT: ****\n");
1138     for(i=0; i< NO_NODES; i++)
1139     {
1140         if(solution_array[i] == 1)
1141             printf("\tNODE %d\n", i+1);
1142     }
1143     printf("\n*** TOTAL COST: %d \n", min_cost);

```

```

1155 )
1156 /*.....
1157 * FIND_COST():
1158 *
1159 * Processing: Called by the branch bound function, this function computes the
1160 * new costs associated with file replication. This is done by parsing
1161 * the query and update cost tables - keeping the proposed replicaion
1162 * of a partition in mind (i.e., if the file were placed on 2 nodes
1163 * queries are now satisfied (or free) at both these nodes, the other
1164 * queries will access the file copy from the node representing the
1165 * least cost. Updates on the file copy from the other hand would be higher as each copy
1166 * of the file must be updated).
1167 *
1168 * Input: The node being analyzed for replication
1169 *
1170 * Functions
1171 * Called: (None)
1172 *
1173 * Output:
1174 * Return: Returns the cost associated with the proposed replication
1175 * .....
1176 */
1177
1178 find_cost(new_node)
1179 {
1180     int calc_qry = 0, calc_upd = 0, tot_cost = 0;
1181     int i, j;
1182
1183     /* DETERMINE QUERY COST WITH NEW FILE PLACEMENT */
1184     for(i=0; i< NO_NODES; i++)
1185     {
1186         if(( i != new_node) && (solution_array[i] != 1))
1187         {
1188             calc_qry = qry_cost[i][new_node];
1189             for(j=0; j< NO_NODES; j++)
1190             {
1191                 if( i != j )
1192                 {
1193                     if(qry_cost[i][j] < calc_qry)
1194                     {
1195                         calc_qry = qry_cost[i][j];
1196                     }
1197                 }
1198             }
1199         }
1200     }
1201 }
1202

```

```

1203     )
1204   } End of "for" Loop */
1205
1206   /* DETERMINE UPDATE COSTS WITH NEW FILE PLACEMENT */
1207
1208   for(i=0; i< NO_NODES; i++)
1209   {
1210     if(( !new_node) || (solution_array[i] == i))
1211     {
1212       for(j=0; j< NO_NODES; j++)
1213         calc_upd = calc_upd + upd_cost[j][i];
1214     }
1215   }
1216
1217   /* TOTAL COST TO RETURN IS QUERY + UPDATE */
1218   tot_cost = calc_qry + calc_upd;
1219
1220   return(tot_cost);
1221 }
1222
1223 /******
1224 * REBUILD_COST():
1225 * Processing: Rebuilds the query and update cost tables based upon a percent
1226 *               increase/decrease in communication costs (called if user desired
1227 *               sensitivity analysis to be performed; in which case the user supplies
1228 *               the percentage of increase/decrease desired)
1229 * Input:        A sign to indicate an increase or decrease is requested, and the
1230 *               percentage to be associated with this change.
1231 * Functions
1232 * Called:       (None)
1233 * Output/
1234 * Return:       (None)
1235 * *****
1236 rebuild_cost(sign,percent)
1237 int sign,percent;
1238 {
1239   int i,j;
1240   if(sign == '+')
1241     svmin_cost = svmin_cost + ((svmin_cost * percent)/100);
1242   else
1243     svmin_cost = svmin_cost - ((svmin_cost * percent)/100);
1244 }
1245
1246
1247
1248
1249
1250

```



```

1251 for(i=0; i<ND_NODES; i++)
1252 {
1253     for(j=0; j<ND_NODES; j++)
1254     {
1255         if(sign == '+')
1256         {
1257             qrv_cost[i][j] = qrv_cost[i][j] + ((qrv_cost[i][j] * percent)/100);
1258             upd_cost[i][j] = upd_cost[i][j] + ((upd_cost[i][j] * percent)/100);
1259         }
1260         else
1261         {
1262             qrv_cost[i][j] = qrv_cost[i][j] - ((qrv_cost[i][j] * percent) / 100);
1263             upd_cost[i][j] = upd_cost[i][j] - ((upd_cost[i][j] * percent)/100);
1264         }
1265     }
1266 }
1267 }
1268 }
1269 }
1270 /*****
1271  * PR_RELTAB():
1272  * Processing: Prints out the contents of the relation table, which indicates
1273  * the attribute usage per node/document for any given relation.
1274  *
1275  * Input: Relation number
1276  *
1277  * Functions
1278  * Called: None
1279  *
1280  * Output/
1281  * Return: print display to user
1282  *
1283  *****/
1284 pr_reltbl(relno)
1285 int relno;
1286 {
1287     int i, j, k;
1288     printf("RELNO =%d, NUMBER OF RELATIONS = %d\n", relno, attrnum[relno]);
1289     for(i=0; i<ND_ENTS; i++)
1290     {
1291         printf("\nNODE: %d ORG: %s QDC: %s\n", reltable[i].nodenum, reltable[i].orgname,
1292             reltable[i].docname);
1293         for(j=0; j<attrnum[relno]+1; j++)
1294         {
1295             printf("%s: ", reltable[relno].attributes[j].attr);
1296             if(reltable[i].freq[j].part_num==ND_NODES + 1)
1297                 printf(" Partition: *KEY* ");
1298         }
1299     }

```

```

1299         printf(" Partition:\t%d",reliable[i].freq[j].part_freq);
1300
1301         if(reliable[i].freq[j].owner == TRUE)
1302             printf("\t%d\n",reliable[i].freq[j].frequency);
1303         else
1304             printf("\t%d\n",reliable[i].freq[j].frequency);
1305     }
1306 }
1307 }

```

Dist_Data.h

```

1 #define MAX_ENTRY 200 /* Max. no. of org/doc entries */
2 #define DDC_SZ 30 /* Max. no. of chars in document name */
3 #define DRG_SZ 30 /* Max. no. of chars in organization name */
4 #define ATR_SZ 20 /* Max. no. of chars in attribute name */
5 #define MAX_ATTRS 300 /* Max. no. of attributes handled */
6 #define MAX_RELS 100 /* Max. no. of relations handled */
7 #define MAX_PARTS 50 /* Max. no. of partitions handled */
8 #define MAX_NODES 50 /* Max. no. of nodes handled */
9 #define MAX_DDCS 200 /* Max. no. of documents handled */
10 #define TRUE 0
11 #define FALSE 1
12
13 /* RELATION TABLE: POPULATED BY NODE, DOCUMENTS and ATTRIBUTES USED */
14
15 struct REL_TABLE
16 {
17     int nodenum;
18     char orgname[DRG_SZ];
19     char docname[DDC_SZ];
20     int use_rel;
21     int total_query;
22     int total_update;
23     struct FREQ
24     {
25         int key;
26         int owner;
27         int part_num;
28         int frequency;
29     } freq[MAX_ATTRS];
30 } reltable[MAX_ENTRY];
31
32 /* BERNTABLE: STORES BERNT2 RELATIONS */
33
34 struct BERNTABLE
35 {
36     struct ATTRIBUTES
37     {
38         char attr[ATR_SZ];
39         int key;
40         attributes[MAX_ATTRS];
41     } berntable[MAX_RELS];

```

```

42 /* TRANSACTION TABLE: HOLDS TOTAL NO. OF QUERY/UPDATE ACCESSES PER NODE */
43
44 struct TRANS
45 {
46     int t_qry;
47     int t_upd;
48     int t_obj[max_nodes];
49 }trans_tbl[max_nodes];
50
51 /* TRANSACTION COST TABLES (update and query) */
52
53 int qry_cost[max_nodes][max_nodes];
54 int upd_cost[max_nodes][max_nodes];
55
56 /* NETWORK COST TABLE */
57
58 int costable[max_nodes][max_nodes];
59
60 /* SOLUTION ARRAY FOR BRANCH & BOUND ALGORITHM */
61
62 int solution_array[max_nodes];
63
64 /* GLOBAL 'HOLDER' FOR MIN_COST PARTITION/FILE PLACEMENT */
65
66 int min_cost, smin_cost, svnode;
67
68 /* NO. OF ATTRIBUTES FOR EACH BERN2 RELATION */
69
70 int attrnum[max_rels];
71
72 /* NUMBER OF PARTITIONS CREATED FROM BERN2 RELATION(s) */
73
74 unsigned int parts_created[max_nodes];
75

```

DATA ALLOCATION IN A
DISTRIBUTED DATABASE ENVIRONMENT

by

Kimberley Ann Johnson

B. S., Western Illinois University, 1979

AN ABSTRACT OF A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Computer Science

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1988

Distributed information systems are systems which locally or geographically distribute elements within a computing system. In this environment the database designer is faced with many new design problems, one of the most critical being the distribution of data which most accurately reflects the processing needs of the organization.

The main focus of this report, is in the design of an interactive system that automates the partitioning of files and their placement onto a distributed network. A review of the relevant literature is first presented, followed by an overview of the system's design. Finally, an extensive example is provided to demonstrate the functionality and use of the interactive system.